



# Algorithmic construction of optimal designs on compact sets for concave and differentiable criteria

Luc Pronzato, Anatoly A. Zhigljavsky

## ► To cite this version:

Luc Pronzato, Anatoly A. Zhigljavsky. Algorithmic construction of optimal designs on compact sets for concave and differentiable criteria. *Journal of Statistical Planning and Inference*, 2014, 154, pp.141-155. 10.1016/j.jspi.2014.04.005 . hal-01001706

**HAL Id: hal-01001706**

**<https://hal.science/hal-01001706>**

Submitted on 4 Jun 2014

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Algorithmic construction of optimal designs on compact sets for concave and differentiable criteria\*

L. PRONZATO<sup>1,2</sup> and A. ZHIGLJAVSKY<sup>3</sup>

<sup>1</sup> Author for correspondence

<sup>2</sup> Laboratoire I3S, CNRS/Université de Nice-Sophia Antipolis  
Bât. Euclide, Les Algorithmes, 2000 route des lucioles, BP 121  
06903 Sophia Antipolis cedex, France

<sup>3</sup> School of Mathematics, Cardiff University  
Senghennydd Road, Cardiff, CF24 4YH, UK

`pronzato@i3s.unice.fr`

`ZhigljavskyAA@cf.ac.uk`

February 19, 2014

## Abstract

We consider the problem of construction of optimal experimental designs (approximate theory) on a compact subset  $\mathcal{X}$  of  $\mathbb{R}^d$  with nonempty interior, for a concave and Lipschitz differentiable design criterion  $\phi(\cdot)$  based on the information matrix. The proposed algorithm combines (a) convex optimization for the determination of optimal weights on a support set, (b) sequential updating of this support using local optimization, and (c) finding new support candidates using properties of the directional derivative of  $\phi(\cdot)$ . The algorithm makes use of the compactness of  $\mathcal{X}$  and relies on a finite grid  $\mathcal{X}_\ell \subset \mathcal{X}$  for checking optimality. By exploiting the Lipschitz continuity of the directional derivatives of  $\phi(\cdot)$ , efficiency bounds on  $\mathcal{X}$  are obtained and  $\epsilon$ -optimality on  $\mathcal{X}$  is guaranteed. The effectiveness of the method is illustrated on a series of examples.

**keywords** Approximate design; optimum design; construction of optimal designs, global optimization

**MSC** 62K05; 90C30; 90C26

---

\*This work was initiated while the authors were invited at the Isaac Newton Institute for Mathematical Sciences, Cambridge, UK; the support of the INI and of CNRS is gratefully acknowledged.

# 1 Introduction and motivation

A design measure  $\xi$  on a finite set  $\mathcal{X} \subset \mathbb{R}^d$  with  $\ell$  elements is characterized by the  $\ell$ -dimensional vector of weights  $W$  (nonnegative and summing to one) allocated to the  $\ell$  elements of  $\mathcal{X}$ . The determination of an optimal measure  $\xi^*$  which maximizes a concave differentiable criterion  $\phi(\cdot)$  then forms a finite-dimensional convex problem for which many optimization algorithms are available, see, *e.g.*, Hiriart-Urruty and Lemaréchal (1993); den Hertog (1994); Nesterov and Nemirovskii (1994); Ben-Tal and Nemirovski (2001); Boyd and Vandenberghe (2004); Nesterov (2004) for recent developments on convex programming. In particular, the cutting plane method of Kelley (1960) is considered by Sibson and Kenny (1975), and a variant of it (with a modified version of the Equivalence Theorem) by Gribik and Kortanek (1977); see also Pronzato and Pázman (2013, Chap. 9). However, design problems are usually such that (a) the cardinality  $\ell$  of the design space, which determines the dimension of the optimization problem to be solved, is large, and (b) there always exists an optimal measure  $\xi^*$  with a few support points only, *i.e.*, such that only a few components of  $W$  are positive (say  $m$ , with  $m \ll \ell$ ). These particularities have motivated the development of specific methods which happen to be competitive compared to general-purpose convex-programming algorithms. One may distinguish three main families, see Pronzato and Pázman (2013, Chap. 9) for a recent survey.

- (i) *Vertex-direction* methods only increase one component of  $W$  at each iteration, all other components being multiplied by the same factor smaller than one; see Wynn (1970); Fedorov (1972) and also Frank and Wolfe (1956) for a method originally proposed in a more general context. A significant improvement can be achieved by decreasing some individual components of  $W$  only, which allows us to set some components of  $W$  to zero, *i.e.*, to remove support points from a poorly chosen initial design and sometimes to exchange two components of  $W$  (vertex-exchange methods); see Atwood (1973); St. John and Draper (1975); Böhning (1985, 1986); Molchanov and Zuyev (2001, 2002).
- (ii) In *gradient* methods the gradient direction projected on the set of design measures is used as the direction of movement at each iteration; see Wu (1978a,b), and Atwood (1976) for an extension to a Newton-type method.
- (iii) At each iteration of a *multiplicative* method, each component of  $W$  is multiplied by a suitably chosen positive function; see, *e.g.*, Titterton (1976); Silvey et al. (1978); Torsney (1983); Fellman (1989); Dette et al. (2008); Yu (2010a,b), and Torsney (2009) for a historical review.

In multiplicative methods all initial weights must be positive, those which should be zero at the optimum decrease continuously along iterations but typically stay strictly positive, *i.e.*,

never achieve exactly zero. The convergence is inevitably slow close to the optimum. The same phenomenon occurs for vertex-direction methods, unless the decrease of some particular individual components of  $W$  is allowed at some iterations, so that poor initial support points can be removed. Even in this case, if  $\ell$  is large many iterations are required to identify the components of  $W$  which should be positive at the optimum. Gradient-type methods are also not efficient when  $\ell$  is large.

Several authors tried combinations of different methods to make use of their respective advantages. A very sensible algorithm is proposed in (Wu, 1978a,b); it combines a gradient method (always working in a small dimensional subspace) with a vertex-direction method (which allows a suitable updating of this subspace). A mixture of multiplicative and vertex-exchange algorithms is proposed in (Yu, 2011) for  $D$ -optimum design; it includes a nearest-neighbor exchange strategy which helps apportioning weights between adjacent points in  $\mathcal{X}$  and has the property that poor support points are quickly removed from the support of the initial measure. Attractive performance is reported. All the methods above, however, are restricted to the case where the design space  $\mathcal{X}$  is finite with  $l$  being not too large and, for some of them, to particular design criteria ( $D$ -optimality, for instance).

When one is interested in the determination of an optimal design on a compact subset  $\mathcal{X}$  of  $\mathbb{R}^d$  with nonempty interior, the usual practice consists in discretizing  $\mathcal{X}$  into a finite set  $\mathcal{X}_\ell$  with  $\ell$  elements and applying one of the methods above. When a precise solution is required, then  $\ell$  is necessarily very large (in some cases,  $\ell = 10^6$  should be considered as a small number) and none of the methods above is efficient. Refining iteratively a finite grid contained in  $\mathcal{X}$  is a possible option, see Wu (1978a), but the search for the optimal design is still performed in a discrete set.

In contrast, the algorithm we propose makes use of the compactness of  $\mathcal{X}$  and looks for the optimal support in the whole set  $\mathcal{X}$ . A finite grid  $\mathcal{X}_\ell \subset \mathcal{X}$  is only used to check optimality on  $\mathcal{X}_\ell$  and, using the Lipschitz continuity of directional derivatives of  $\phi(\cdot)$ , to construct an efficiency bound on  $\mathcal{X}$ . A key ingredient in the algorithm is the separation between the determination of the support points of  $\xi^*$  (knowing that there are at most  $m$  of them) from the determination of the associated weights (an  $m$ -dimensional convex problem).

The determination of the support of  $\xi^*$  is a non-convex problem, usually multimodal, for which the straightforward application of a global search method (such as the simulated annealing or one of the genetic algorithms) cannot be recommended. Indeed, these heuristic methods do not use the crucial information about the objective function provided by the Lipschitz constants, derivatives and convexity of the weight-optimization problem, and most of them do not provide any indication of the closeness of the returned solution to an optimum, global or local. On the other hand, by using properties of the directional derivative of  $\phi(\cdot)$ , we can easily locate good

candidates for the support of  $\xi^*$  and thereby we do not need to perform a global search in the  $m \times d$  dimensional space  $\mathcal{X}^d$ . The aim of this paper is to show how these properties can be combined with algorithms for convex optimization which are classically used in design for the determination of optimal weights to yield an efficient method of construction of optimal designs on compact sets.

The paper is organized as follows. Section 2 defines the problem and introduces the notation. Section 3 presents the algorithm and proves its convergence. A few illustrative examples are given in Sect. 4 where the results obtained with the proposed algorithm are discussed. Section 5 concludes and indicates some remaining open issues. A few technical aspects are collected in Appendix.

## 2 Notation and problem statement

Let the design space  $\mathcal{X}$  be a compact subset of  $\mathbb{R}^d$  with nonempty interior (typically  $\mathcal{X} = [-1, 1]^d$ ) and denote by  $\Xi(\mathcal{X})$  the set of probability measures on  $\mathcal{X}$ . Any element  $\xi \in \Xi(\mathcal{X})$  will be called design measure, or shortly design.

Consider a model (the main example being a regression model) defined on  $\mathcal{X}$  with  $p$  unknown parameters  $\theta \in \mathbb{R}^p$ ,  $p \geq 2$ . Denote by  $\mathbb{M}^{\geq}$  and  $\mathbb{M}^{>}$  the sets of  $p \times p$  symmetric non-negative and positive definite matrices respectively. Let

$$M(\xi) = \int_{\mathcal{X}} \mu(x) \xi(dx) \quad (1)$$

be the information matrix related to the estimation of  $\theta$ , with  $\mu(x) \in \mathbb{M}^{\geq}$ . The matrix  $\mu(x)$  is the elementary information matrix for the estimation of  $\theta$  associated with one single observation at  $x$  and may have rank larger than one (think of Bayesian optimal design or multivariate regression, see, *e.g.*, Fedorov (1972, Chap. 5), Harman and Trnovská (2009)). In the case of simple linear regression model with observations  $y_j = \theta^T f(x_j) + \varepsilon_j$ , where the  $\varepsilon_j$  denote i.i.d. random errors with zero mean, we have  $\mu(x) = f(x)f^T(x)$ . We suppose that  $\mu(\cdot)$  is continuous on  $\mathcal{X}$  (this assumption could be weakened, but it would make the presentation unduly complicated).

Although  $\mu(\cdot)$ , and thus  $M(\xi)$ , may depend on the model parameters  $\theta$  in a nonlinear situation, this dependence will be omitted and, when it occurs, we only consider locally optimum designs and evaluate  $M(\cdot)$  at a nominal value for  $\theta$ . The set  $\mathcal{M}_{\mathcal{X}} = \{M(\xi) : \xi \in \Xi(\mathcal{X})\}$  is a convex and compact subset of a linear space of dimension  $p(p+1)/2 + 1$ ; from Caratheodory's Theorem, see, *e.g.*, Silvey (1980, p. 72), any element of  $\mathcal{M}_{\mathcal{X}}$  can be written as  $M(\xi_{\text{discr}})$ , where  $\xi_{\text{discr}}$  is a discrete measure supported on a finite set with  $m \leq p(p+1)/2 + 1$  elements.

We are interested in constructing a  $\phi$ -optimal design on  $\mathcal{X}$ ; that is, a design  $\xi^* \in \Xi(\mathcal{X})$

which maximizes the functional

$$\phi(\xi) = \Phi[M(\xi)]$$

with respect to  $\xi \in \Xi(\mathcal{X})$ . The function  $\Phi : \mathbb{M}^{\geq} \rightarrow \mathbb{R}$  is normally assumed to be (see, *e.g.*, Pukelsheim (1993, Chap. 5)):

- (i) concave, *i.e.*,  $\Phi[(1 - \alpha)M_1 + \alpha M_2] \geq (1 - \alpha)\Phi(M_1) + \alpha\Phi(M_2)$  for all  $M_1, M_2$  in  $\mathbb{M}^{\geq}$  and all  $\alpha$  in  $[0, 1]$ ;
- (ii) proper, *i.e.*, the set  $\{M \in \mathbb{M}^{\geq} : \Phi(M) > -\infty\}$  is nonempty and  $\Phi(M) < \infty$  for all  $M \in \mathbb{M}^{\geq}$ ;
- (iii) isotonic (monotonic for Loewner ordering), *i.e.*,  $\Phi(M_1) \geq \Phi(M_2)$  for all  $M_1, M_2$  in  $\mathbb{M}^{\geq}$  such that  $M_1 - M_2$  is non-negative definite.

We suppose that a nonsingular design measure  $\xi_0$  exists in  $\Xi(\mathcal{X})$  (*i.e.*, such that  $M(\xi_0) \in \mathbb{M}^>$ ). We shall denote by  $\Phi^+(\cdot)$  the positively homogeneous version of  $\Phi(\cdot)$ , which satisfies  $\Phi^+(I_p) = 1$ , with  $I_p$  the  $p$ -dimensional identity matrix, and  $\Phi^+(aM) = a\Phi^+(M)$  for all  $M$  in  $\mathbb{M}^{\geq}$  and all  $a$  positive. We also suppose that the design criterion  $\Phi(\cdot)$  is global, that is, its positively homogeneous version  $\Phi^+(\cdot)$  satisfies  $\Phi^+(M) > 0$  if and only if  $M \in \mathbb{M}^>$ . We can write  $\Phi^+(M) = \psi[\Phi(M)]$ , with  $\psi(\cdot)$  the inverse function of the transformation  $a \in \mathbb{R}^+ \mapsto \Phi(aI_p)$ . Note that  $\psi(\cdot)$  is increasing. We suppose that  $\Phi^+(\cdot)$  is concave. Standard examples are  $\Phi(M) = \log \det(M)$  with  $\Phi^+(M) = \det^{1/p}(M)$ ,  $\psi(t) = \exp(t/p)$ , and  $\Phi(M) = -\text{trace}(M^{-1})$  with  $\Phi^+(M) = [(1/p)\text{trace}(M^{-1})]^{-1}$  and  $\psi(t) = -p/t$ ; see, *e.g.*, Pukelsheim (1993, Chap. 5). The  $\phi$ -efficiency of a design measure  $\xi \in \Xi(\mathcal{X})$  is defined as

$$\mathcal{E}_{\phi}(\xi) = \frac{\phi^+(\xi)}{\max_{\xi \in \Xi(\mathcal{X})} \phi^+(\xi)}, \quad (2)$$

where  $\phi^+(\xi) = \Phi^+[M(\xi)]$ .

We assume that  $\Phi(\cdot) : \mathbb{M}^{\geq} \rightarrow \mathbb{R}$  is differentiable on  $\mathbb{M}^>$  with Lipschitz continuous gradient; that is,

$$\|\nabla\Phi(M_1) - \nabla\Phi(M_2)\| \leq L(c)\|M_1 - M_2\| \quad (3)$$

for all  $M_1, M_2$  in the set

$$\mathcal{M}(c) = \{M \in \mathbb{M}^> : \Phi(M) \geq c\},$$

for some constant  $L(c)$ , where  $\nabla\Phi(M)$  denotes the gradient of  $\Phi(\cdot)$  at  $M$ , the norm  $\|\cdot\|$  is the usual  $L_2$ -norm (for any real matrix  $M$ ,  $\|M\| = \text{trace}^{1/2}(M^{\top}M)$ ) and  $c$  is any number between  $\inf_{\xi \in \Xi(\mathcal{X})} \phi(\xi)$  and  $\max_{\xi \in \Xi(\mathcal{X})} \phi(\xi)$ . We also suppose that

$$\|\nabla\Phi(M)\| \leq B(c) \text{ for all } M \in \mathcal{M}(c) \quad (4)$$

for some constant  $B(c)$ . Note that  $L(c)$  and  $B(c)$  will only be used in the proofs of Theorems 2 and 3 and that their knowledge is not required for the construction of optimal designs with the algorithms of Sect. 3.

Isotonicity implies that any optimal information matrix  $M(\xi^*)$  lies on the boundary of the set  $\mathcal{M}_{\mathcal{X}}$ , and therefore there always exists an optimal design measure with no more than  $m^*(p) = p(p+1)/2$  support points. We can thus restrict our attention to such designs, *i.e.*, to designs  $\xi$  defined by their support points  $\{x^{(1)}, \dots, x^{(m)}\}$ ,  $m \leq m^*(p)$ , and associated weights  $w_i = \xi(\{x^{(i)}\}) \geq 0$  which satisfy  $\sum_{i=1}^m w_i = 1$ . The notation

$$\xi = \begin{Bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ w_1 & w_2 & \dots & w_m \end{Bmatrix}$$

is standard. One should notice that the bound  $m^*(p)$  is often pessimistic since many situations exist where optimal designs are concentrated on a much smaller number of support points or are even saturated, *i.e.*, are supported on  $p$  points only; see Yang and Stufken (2009); Yang (2010) and the paper (Dette and Melas, 2011) which makes use of results in (Karlin and Studden, 1966b) on Chebyshev systems and moment spaces.

For any  $n \geq 1$  and any design measure  $\xi$  on  $X = \{x^{(1)}, \dots, x^{(n)}\} \subset \mathcal{X}^n$ , we shall denote by  $W = W(\xi) = (w_1, \dots, w_n)^\top \in \mathcal{P}_n$  the vector formed by the weights  $w_i = \{W\}_i = \xi(\{x^{(i)}\})$ , with  $\mathcal{P}_n$  the probability simplex

$$\mathcal{P}_n = \{W = (w_1, \dots, w_n)^\top \in \mathbb{R}^n : w_i \geq 0, \sum_{i=1}^n w_i = 1\}. \quad (5)$$

Some components  $w_i$  of  $W$  may equal zero, and we shall denote by  $W^+(\xi)$  the vector formed by strictly positive weights and by  $S = S(\xi)$  the support of  $\xi$ , defined by the corresponding design points. With a slight abuse of notation, we shall consider  $X$  and  $S$  sometimes as sets with respectively  $n$  and  $m \leq n$  elements, and sometimes as vectors in  $\mathbb{R}^{n \times d}$  and  $\mathbb{R}^{m \times d}$ . Also, we shall respectively denote by  $\phi(S|W^+)$  and  $\phi(W^+|S) = \phi(W|X)$  the value of the design criterion  $\phi(\xi)$  when considered as a function of  $S$  with  $W^+$  fixed and as a function of  $W^+$  with  $S$  fixed, and by  $\nabla \phi(W^+|S)$  the usual gradient  $\partial \phi(W^+|S) / \partial W^+$  (a  $m$ -dimensional vector). For any  $\nu$  and  $\xi$  in  $\Xi(\mathcal{X})$  such that  $\phi(\xi) > -\infty$ , we can compute the directional derivative of  $\phi(\cdot)$  at  $\xi$  in the direction  $\nu$ ,

$$F(\xi; \nu) = \lim_{\alpha \rightarrow 0^+} \frac{\phi[(1-\alpha)\xi + \alpha\nu] - \phi(\xi)}{\alpha} = \text{trace}\{[M(\nu) - M(\xi)]\nabla \Phi[M(\xi)]\}$$

(with  $F(\xi; \nu) \in \mathbb{R} \cup \{+\infty\}$ ). When  $\nu$  is the delta-measure  $\delta_x$  at  $x$ , we obtain

$$F(\xi, x) = F(\xi; \delta_x) = \text{trace}\{[\mu(x) - M(\xi)]\nabla \Phi[M(\xi)]\}. \quad (6)$$

Suppose that  $\xi$  and  $\nu$  are finitely supported and let  $X$  denote the union of their supports. Denoting by  $W$  and  $W'$  the vectors of weights which  $\xi$  and  $\nu$  respectively allocate to points in  $X$ , with some components of  $W$  and  $W'$  possibly equal to zero, we obtain

$$F(\xi; \nu) = (W' - W)^\top \nabla \phi(W|X).$$

The concavity and differentiability of  $\Phi(\cdot)$  yield a necessary and sufficient condition for the optimality of a design measure  $\xi^*$ , see, *e.g.*, Silvey (1980), usually named the Equivalence Theorem as a tribute to the work of Kiefer and Wolfowitz (1960); see also Karlin and Studden (1966a).

**Theorem 1** *The following properties are equivalent*

- (i)  $\xi^*$  is  $\phi$ -optimal on  $\mathcal{X}$ , i.e.,  $\xi^*$  maximizes  $\phi(\xi)$  with respect to  $\xi \in \Xi(\mathcal{X})$ ;
- (ii)  $\xi^*$  minimizes  $\max_{x \in \mathcal{X}} F(\xi, x)$  with respect to  $\xi \in \Xi(\mathcal{X})$ ;
- (iii)  $\max_{x \in \mathcal{X}} F(\xi^*, x) = 0$ .

Note that  $F(\xi; \xi) = 0$  for all  $\xi \in \Xi(\mathcal{X})$  and that  $F(\xi; \nu) = \int_{\mathcal{X}} F(\xi, x) \nu(dx)$  for all  $\xi$  and  $\nu \in \Xi(\mathcal{X})$ , so that (iii) implies that  $\xi^*$ -almost everywhere we have  $F(\xi^*, x) = 0$ ; that is,  $F(\xi^*, x) = 0$  on the support of  $\xi^*$ .

We shall say that  $\xi^\epsilon$  is  $\epsilon$ -optimal for  $\phi(\cdot)$  on  $\mathcal{X}$  if and only if

$$\max_{x \in X} F(\xi^\epsilon, x) < \epsilon.$$

The concavity of  $\phi(\cdot)$  then implies that  $\phi(\xi^\epsilon) > \max_{\xi \in \Xi(X)} \phi(\xi) - \epsilon$ .

The directional derivative of the positively homogeneous criterion  $\phi^+(\cdot) = \psi[\phi(\cdot)]$  at  $\xi$  in the direction  $\nu$  is  $F_{\phi^+}(\xi; \nu) = \psi'[\phi(\xi)]F(\xi; \nu)$ . We denote  $F_{\phi^+}(\xi, x) = F_{\phi^+}(\xi; \delta_x)$  for all  $x \in \mathcal{X}$  and all  $\xi \in \Xi(\mathcal{X})$ . Due to the concavity of  $\phi^+(\cdot)$ , the  $\epsilon$ -optimality of  $\xi^\epsilon$  for  $\phi(\cdot)$  implies that  $\phi^+(\xi^\epsilon) > \max_{\xi \in \Xi(X)} \phi^+(\xi) - \epsilon \psi'[\phi(\xi^\epsilon)]$ ; this yields the efficiency bound

$$\mathcal{E}_\phi(\xi^\epsilon) > 1 - \frac{\epsilon \psi'[\phi(\xi^\epsilon)]}{\max_{\xi \in \Xi(X)} \phi^+(\xi)} \geq \underline{\mathcal{E}}_\phi(\xi^\epsilon) = 1 - \frac{\epsilon \psi'[\phi(\xi^\epsilon)]}{\phi^+(\xi^\epsilon)}, \quad (7)$$

where the  $\phi$ -efficiency  $\mathcal{E}_\phi(\cdot)$  is defined in (2).

### 3 Algorithms

We denote by  $A_0$  a convex optimization algorithm for the determination of optimal weights. When initialized at an arbitrary nonsingular design on  $X$ , for any  $\epsilon > 0$ , the algorithm  $A_0$  is assumed to return an  $\epsilon$ -optimal design in a finite number of iterations; we shall denote this  $\epsilon$ -optimal design by  $\xi = A_0[X, \epsilon]$ . In the examples of Sect. 4, we use a combination of the projected-gradient and vertex-exchange methods, following ideas similar to those in (Wu, 1978a), but other methods could be used as well, see a discussion in Sect. 1.



### Algorithm $A_0$

**Input:** Discrete set  $X = \{x^{(1)}, \dots, x^{(m)}\} \in \mathcal{X}^m$  (such that  $M(\nu_0)$  has full rank, with  $\nu_0$  the uniform measure on  $X$ ) and  $\epsilon > 0$ .

**Output:** Design  $\xi = A_0[X, \epsilon]$  such that  $\max_{x \in X} F(\xi, x) < \epsilon$ .

### 3.1 Construction of the main algorithm

We shall construct an algorithm for the maximization of  $\phi(\cdot)$  on the set of design measures on the compact set  $\mathcal{X}$ . We denote by  $\xi_k$  the design measure obtained at iteration  $k$  (always with finite support), by  $S_k = S(\xi_k)$  its support and by  $W_k^+$  the vector of associated weights, and write  $\phi_k = \phi(\xi_k)$ .

The construction is decomposed into three steps. First, we present a prototype algorithm  $A_1$ ; at each iteration  $k$  of this algorithm we determine an  $(\epsilon/2)$ -optimal design  $\xi_k = A_0[X_k, \epsilon/2]$  for an increasing sequence of sets  $X_k$ , where  $X_k$  has one element more than  $X_{k-1}$ . The main difference between algorithm  $A_1$  and the vertex-direction methods discussed in Sect. 1 is step 1, where the weights at the current support are optimized. Although closely related, the algorithm also differs from a direct application of the cutting-plane method as considered by Sibson and Kenny (1975); Gribik and Kortanek (1977).

Algorithm  $A_2$  is very similar to  $A_1$  but in  $A_2$  the number of points in the sets  $X_k$  is bounded.

In both  $A_1$  and  $A_2$  we assume that  $\max_{x \in \mathcal{X}} F(\xi_k, x)$  can be determined with arbitrary precision  $\epsilon$ . In practice, we usually perform the search for the maximum of  $F(\xi_k, x)$  over a finite test-set  $\mathcal{X}_\ell$ . Using Lipschitz continuity arguments and (7), we can then guarantee some efficiency bound over  $\mathcal{X}$  when  $\mathcal{X}_\ell \subset \mathcal{X}$  is a fine enough grid, see Sect. 3.3. However, the computational effort required becomes too high when  $\ell$  gets large if we search for the maximum of  $F(\xi_k, x)$  over the whole grid  $\mathcal{X}_\ell$  at each iteration. In algorithm  $A_3$ , we use a two-level strategy where the maximization of  $F(\xi_k, \cdot)$  is performed over a coarse test-set  $\mathcal{T}_k$  which is progressively enriched by points taken from  $\mathcal{X}_\ell$ . Moreover, the maximization over  $\mathcal{T}_k$  is complemented by a local search over  $\mathcal{X}$ .

The addition of a local maximization of  $\phi(S|W_k^+)$  with respect to the support  $S$  of  $\xi_k$  at each iteration preserves the convergence of  $\phi_k$  to the optimum  $\phi^* = \phi(\xi^*)$  and yields our final algorithm  $A_4$ , which also incorporates a merging strategy which avoids having clusters of points (that is, groups of points positioned closely) in  $\mathcal{T}_k$  and  $X_k$ . This is important as the presence of such clusters slows down the convergence of algorithms.

**Algorithm  $A_1$** 

**Step 0)** Choose  $X_0 = \{x^{(1)}, \dots, x^{(m)}\} \in \mathcal{X}^m$  such that  $M(\nu_0)$  has full rank, with  $\nu_0$  the uniform measure on  $X_0$ ; choose  $\epsilon > 0$ , set  $k = 0$ .

**Step 1)** Compute  $\xi_k = A_0[X_k, \epsilon/2]$ .

**Step 2)** Find  $x_k^* = \arg \max_{x \in \mathcal{X}} F(\xi_k, x)$ .

**Step 3)** If  $F(\xi_k, x_k^*) < \epsilon$ , stop; otherwise, set  $X_{k+1} = \{X_k, x_k^*\}$ ,  $k \leftarrow k + 1$ , go to step 1.

Note that we can always use  $m \leq p(p+1)/2$  at step 0 above. The next theorem establishes an important property of algorithm  $A_1$ .

**Theorem 2** *Algorithm  $A_1$  stops after a finite number of iterations. The design measure  $\xi_k$  obtained at step 3 is then  $\epsilon$ -optimal, in the sense  $\phi_k = \phi(\xi_k) > \phi^* - \epsilon$ .*

*Proof.* Denote by  $\phi_k^*$  the maximum value of  $\phi(\cdot)$  on the set of probability measures supported on  $X_k$ . From the  $(\epsilon/2)$ -optimality of the design  $\xi_k$  returned by algorithm  $A_0$ , we have  $\phi_0 = \phi(\xi_0) \geq \phi_0^* - \epsilon/2$ . Since  $X_0 \subset X_k$  for any  $k \geq 1$ ,  $\phi_k^* \geq \phi_0^*$  and  $\phi_k \geq \phi_k^* - \epsilon/2 \geq \phi_0^* - \epsilon/2$  for all  $k$ .

Now, for all  $k$  and all  $x, x' \in \mathcal{X}$  we have

$$\begin{aligned} |F(\xi_k, x') - F(\xi_k, x)| &= |\text{trace}\{[\mu(x') - \mu(x)] \nabla \Phi[M(\xi_k)]\}| \\ &\leq \|\mu(x') - \mu(x)\| \|\nabla \Phi[M(\xi_k)]\| \\ &\leq \|\mu(x') - \mu(x)\| B(\phi_0^* - \epsilon/2), \end{aligned}$$

see (4) and (6). Since  $\mathcal{X}$  is compact and  $\mu(x)$  is continuous, there exists  $\rho > 0$  such that for any  $x$  and  $x'$  in  $\mathcal{X}$ ,  $\|x' - x\| < \rho$  implies that  $|F(\xi_k, x') - F(\xi_k, x)| < \epsilon/2$  for all  $k$ . Since  $\xi_k = A_0[X_k, \epsilon/2]$ ,  $F(\xi_k, x^{(i)}) < \epsilon/2$  for all  $x^{(i)}$  in  $X_k$ , so that any  $x$  in a ball  $\mathcal{B}(x^{(i)}, \rho)$  with centre  $x^{(i)}$  and radius  $\rho$  is such that  $F(\xi_k, x) < \epsilon$ . This implies that  $x_k^*$  constructed at step 2 is at distance larger than  $\rho$  from all points in  $X_k$  when algorithm  $A_1$  does not stop at step 3. In view of the fact that  $\mathcal{X}$  is compact,  $A_1$  necessarily stops after a finite number of iterations.

The concavity of  $\phi(\cdot)$  implies that  $\phi^* \leq \phi_k + \max_{x \in \mathcal{X}} F(\xi_k, x)$  and thus  $\phi^* < \phi_k + \epsilon$  when the algorithm stops. ■

**Remark 1** *Note that the algorithm stops in less than  $k^*$  iterations, with*

$$k^* = \max\{k : \text{there exists } (x_1, \dots, x_k) \in \mathcal{X}^k \text{ such that } \min_{i \neq j} \|x_i - x_j\| \geq \rho\}.$$

*An upper bound on  $k^*$  (very pessimistic) can easily be constructed when  $\mathcal{X}$  is an hypercube in  $\mathbb{R}^d$ , say  $[-1, 1]^d$ , using packing radius. Indeed, the existence of  $(x_1, \dots, x_k)$  in the definition*

of  $k^*$  implies that  $k$  non-intersecting balls of radius  $R = \rho/2$  can be packed in the hypercube  $[-(1 + \rho/2), (1 + \rho/2)]^d$ , so that  $k^* \leq (2 + 4/\rho)^d/V_d$ , with  $V_d = \pi^{d/2}/\Gamma(d/2 + 1)$  the volume of the  $d$ -dimensional unit ball.

In algorithm  $A_1$ , the size of  $X_k$  increases by one at each iteration (see step 3), so that the dimension of the optimization problem to be solved by  $A_0$  at step 1 is also increasing. In order to facilitate the work of  $A_0$  we shall consider another construction for  $X_k$ .

The support  $S_k = S(\xi_k)$  of the design  $\xi_k$  determined at step 1 may be strictly included in  $X_k$ . It is then tempting to replace step 3 of  $A_1$  by

**Step 3')** If  $F(\xi_k, x_k^*) < \epsilon$  stop; otherwise set  $X_{k+1} = \{S_k, x_k^*\}$ ,  $k \leftarrow k + 1$ , go to step 1.

However, the arguments used in the proof of Theorem 2 for the non-existence of cluster points for the sequence  $\{x_k^*\}$  are no longer valid since  $X_k$  does not contain all previous points  $x_{k'}^*$ ,  $k' < k$ . We show in Theorem 3 that convergence is ensured by choosing a suitably decreasing sequence of constants  $\gamma_k$  in  $\xi_k = A_0[X_k, \gamma_k\epsilon]$  at step 1, *e.g.*,  $\gamma_k = 1/(k + 1)$  or  $\gamma_k = \gamma^k$  for some  $\gamma \in (0, 1)$ . The algorithm is then as follows.

**Algorithm  $A_2$**  Let  $\{\gamma_k\}$ ,  $k \geq 0$ , denote a decreasing sequence of numbers in  $(0, 1)$ , with  $\lim_{k \rightarrow \infty} \gamma_k = 0$ .

**Step 0)** Same as in  $A_1$ .

**Step 1)** Compute  $\xi_k = A_0[X_k, \gamma_k\epsilon]$ .

**Step 2)** Same as in  $A_1$ .

**Step 3)** Use step 3' above.

**Theorem 3** *Algorithm  $A_2$  stops after a finite number of iterations and the design measure  $\xi_k$  obtained at step 3 is  $\epsilon$ -optimal.*

*Proof.* We use the same notation as in the proof of Theorem 2. The maximum  $\phi_k^*$  of  $\phi(\cdot)$  in the space of probability measures supported on  $X_k$  is also the maximum of  $\phi(\cdot)$  for a measure supported on  $S_k$ . Since  $S_k \subset X_{k+1}$ ,  $\phi_{k+1}^* \geq \phi_k^*$  and therefore  $\phi_k > \phi_k^* - \gamma_k\epsilon \geq \phi_0^* - \gamma_k\epsilon \geq \phi_0^* - \epsilon$  for all  $k \geq 0$ .

We suppose that  $A_2$  never stops and show that we reach a contradiction. Denote  $\xi_k'(\alpha) = (1 - \alpha)\xi_k + \alpha\delta_{x_k^*}$ , with  $\delta_{x_k^*}$  the delta measure allocating mass 1 at  $x_k^*$ . The function  $\alpha \mapsto \phi[\xi_k'(\alpha)]$  is concave, reaches its maximum for some  $\alpha_k \in [0, 1)$  and equals  $\phi_k$  for  $\alpha = 0$  and  $\alpha = \alpha_k'$ , with  $\alpha_k'$  some number larger than  $\alpha_k$ . For any  $\alpha \in [0, \alpha_k']$  we may write

$$\phi[\xi_k'(\alpha)] = \phi_k + \alpha \text{trace}\{[\mu(x_k^*) - M(\xi_k)]\nabla\Phi[M(\xi_k'(\alpha'))]\}$$

for some  $\alpha' \in [0, \alpha]$ . Using (3) and the fact that  $\phi_k > \phi_0^* - \epsilon$ , we obtain

$$\begin{aligned}
\phi[\xi'_k(\alpha)] &= \phi_k + \alpha \text{trace}\{[\mu(x_k^*) - M(\xi_k)] \nabla \Phi[M(\xi_k)]\} \\
&\quad + \alpha \text{trace}\{[\mu(x_k^*) - M(\xi_k)] [\nabla \Phi[M(\xi'_k(\alpha'))] - \nabla \Phi[M(\xi_k)]]\} \\
&\geq \phi_k + \alpha F(\xi_k, x_k^*) - \alpha \|\mu(x_k^*) - M(\xi_k)\| \|\nabla \Phi[M(\xi'_k(\alpha'))] - \nabla \Phi[M(\xi_k)]\| \\
&\geq \phi_k + \alpha F(\xi_k, x_k^*) - \alpha \alpha' \|\mu(x_k^*) - M(\xi_k)\|^2 L(\phi_0^* - \epsilon) \\
&\geq \phi_k + \alpha \epsilon - \alpha^2 C L(\phi_0^* - \epsilon)
\end{aligned}$$

with  $C = \max_{x, x' \in \mathcal{X}} \|\mu(x') - \mu(x)\|^2$ , where we used the property  $F(\xi_k, x_k^*) \geq \epsilon$ . Therefore,

$$\phi_{k+1}^* \geq \max_{\alpha} \phi[\xi'_k(\alpha)] \geq \phi_k + \frac{\epsilon^2}{4C L(\phi_0^* - \epsilon)}$$

and

$$\phi_{k+1} \geq \phi_{k+1}^* - \gamma_{k+1} \epsilon \geq \phi_k + \frac{\epsilon^2}{4C L(\phi_0^* - \epsilon)} - \gamma_{k+1} \epsilon. \quad (8)$$

This implies that there exists some  $k_0$  such that, for all  $k > k_0$ ,  $\phi_{k+1} > \phi_k + \epsilon^2/[8C L(\phi_0^* - \epsilon)]$ . The sequence  $\{\phi_k\}$  is thus unbounded, which contradicts the assumptions on  $\mathcal{X}$ ,  $\mu(\cdot)$  and  $\phi(\cdot)$ .

Algorithm  $A_2$  thus stops in a finite number of iterations. As in the proof of Theorem 2, the concavity of  $\phi(\cdot)$  implies that  $\xi_k$  is then  $\epsilon$ -optimal.  $\blacksquare$

**Remark 2** An upper bound  $k^*$  on the number of iterations required by  $A_2$  can easily be derived from (8). Indeed,  $k^*$  satisfies  $k^* \epsilon^2 / A - \epsilon \sum_{i=1}^{k^*} \gamma_i \leq \phi^* - \phi_0$ , where we have denoted  $A = 4C L(\phi_0^* - \epsilon)$  and  $\phi^* = \phi(\xi^*)$ . This gives  $k^* \leq (A/\epsilon^2)[\phi^* - \phi_0 + \gamma\epsilon/(1-\gamma)]$  for  $\gamma_i = \gamma^i$  with  $\gamma \in (0, 1)$ , and  $k^* \leq A^2/(4\epsilon^2)[1 + \sqrt{1 + 4(\phi^* - \phi_0)/A}]^2$  for  $\gamma_i = 1/(i+1)$  (where we have used the property  $(1/k) \sum_{i=1}^k \gamma_i < 1/\sqrt{k}$ ).

Step 2 of algorithms  $A_1$  and  $A_2$  may be difficult to implement when  $\mathcal{X}$  has nonempty interior since the function  $x \in \mathcal{X} \mapsto F(\xi_k, x)$  is severely multimodal. We thus discretize  $\mathcal{X}$  into a finite grid  $\mathcal{X}_\ell$  and use the maximizer over this grid as initialization for a local maximization of  $F(\xi_k, \cdot)$  over  $\mathcal{X}$ . The grid must be fine enough (i) to ensure that the maximum of  $F(\xi_k, \cdot)$  is not missed and (ii) to guarantee a reasonable level of optimality over  $\mathcal{X}$ , see Sect. 3.3. However, the computational cost corresponding to the evaluation of  $F(\xi_k, x)$  for all points of  $\mathcal{X}_\ell$  at each iteration would be too high when  $\ell$  is very large. For that reason, algorithm  $A_3$  presented below uses a two-level strategy, with a test-set  $\mathcal{T}_k$  remaining much smaller than  $\mathcal{X}_\ell$ .

We denote by  $x^* = LM[F(\xi, \cdot); x_0]$  the result of a local maximization of  $F(\xi, x)$  with respect to  $x \in \mathcal{X}$  initialized at  $x_0$ , see Sect. 3.2 for possible algorithms when  $\mathcal{X}$  is the hypercube  $[-1, 1]^d$ , or the probability simplex  $\mathcal{P}_d$  given by (5) (mixture experiments). As in  $A_2$ ,  $\{\gamma_k\}$  denotes a decreasing sequence of numbers in  $(0, 1)$ , with  $\lim_{k \rightarrow \infty} \gamma_k = 0$ .

**Algorithm  $A_3$**

- Step 0)** Choose  $X_0 = \{x^{(1)}, \dots, x^{(m)}\} \in \mathcal{X}^m$  such that  $M(\nu_0)$  has full rank, with  $\nu_0$  the uniform measure on  $X_0$ ; choose some arbitrary test-set  $\mathcal{T}_0$  of  $n$  points in  $\mathcal{X}$ ; choose  $\epsilon > 0$ , set  $k = 0$ .
- Step 1)** Compute  $\xi_k = A_0[X_k, \gamma_k \epsilon / 2]$ .
- Step 2)** Find  $x^* = \arg \max_{x \in \mathcal{T}_k} F(\xi_k, x)$  and compute  $x_k^* = LM[F(\xi_k, \cdot); x^*]$ .
- Step 3)** If  $F(\xi_k, x_k^*) \geq \epsilon/2$ , go to step 7; otherwise go to step 4.
- Step 4)** Find  $x^{**} = \arg \max_{x \in \mathcal{X}_\ell} F(\xi_k, x)$ .
- Step 5)** If  $F(\xi_k, x^{**}) < \epsilon/2$ , stop.
- Step 6)** Compute  $x_k^* = LM[F(\xi_k, \cdot); x^{**}]$ .
- Step 7)** Set  $X_{k+1} = \{S_k, x_k^*\}$ ,  $\mathcal{T}_{k+1} = (\mathcal{T}_k, x_k^*)$ ,  $k \leftarrow k + 1$ ; go to step 1.

The test-set  $\mathcal{T}_0$  of step 0 is arbitrary, one may choose for instance  $\mathcal{T}_0 = X_0$  or the  $2^d$  vectors with components  $\pm 1$  when  $\mathcal{X} = [-1, 1]^d$ . The algorithm repeatedly goes through steps 1 to 3 and 7 until  $F(\xi_k, x)$  becomes smaller than  $\epsilon/2$  for all  $x \in \mathcal{T}_k$ ; then the algorithm goes through step 5 which tests the  $(\epsilon/2)$ -optimality of  $\xi_k$  over  $\mathcal{X}_\ell$ . The proof of convergence uses the same arguments as for  $A_2$ : the algorithm stops after a finite number of iterations and returns a design measure  $\xi_k$  that is  $(\epsilon/2)$ -optimal on  $\mathcal{X}_\ell$ . The set  $\mathcal{X}_\ell$  can be chosen such that  $(\epsilon/2)$ -optimality on  $\mathcal{X}_\ell$  implies  $\epsilon$ -optimality on the whole compact set  $\mathcal{X}$ , see Sect. 3.3.

The following two easy improvements can be made to  $A_3$ .

- (i) When  $\min_{x \in S_k} \|x - x_k^*\| < \delta_k$  at step 7, with  $\delta_k$  some small number, one may consider replacing the point in  $S_k$  closest to  $x_k^*$  by  $x_k^*$  and accept the substitution if  $\phi(\xi_k) < \phi(\xi_k')$ , with  $\xi_k'$  the design measure after the substitution. We denote the result of this operation by  $X_{k+1} = EX[\{S_k, x_k^*\}; \delta_k]$ .
- (ii) Step 1 can be complemented by a local maximization of  $\phi(S|W_k^+)$  with respect to  $S \in \mathcal{X}^{m_k}$ , initialized at  $S_k$ , with  $m_k$  the number of support points of  $\xi_k$ ,  $S_k$  their collection and  $W_k^+$  the associated weights. We shall denote the result of this local maximization by  $S_k' = LM[\phi(\cdot|W_k^+); S_k]$ ; see Sect. 3.2 for algorithms.

Since (i) and (ii) above can only increase the value of  $\phi_k$ , the arguments used in the proof of Theorem 3 remains valid and the convergence to an  $(\epsilon/2)$ -optimal measure on  $\mathcal{X}_\ell$  is preserved.

We also add modifications which may decrease the value of  $\phi(\cdot)$  and for which some caution is thus required.

- (iii) If the substitution is rejected in  $X_{k+1} = EX[\{S_k, x_k^*\}; \delta_k]$ , see (i), then the set  $X_{k+1}$  may contain points which are very close, which may complicate the task of  $A_0$  at step 1. We thus consider aggregation of points in  $X_{k+1}$  and replace (sequentially) all pairs of points  $x_i \neq x_j$  such that  $\|x_i - x_j\| < \delta_k$  by  $(x_i + x_j)/2$ , until the set  $X'_{k+1}$  which is obtained is such that  $\min_{x_i \neq x_j \in X'_{k+1}} \|x_i - x_j\| \geq \delta_k$ . We denote this operation by  $X'_{k+1} = AG(X_{k+1}, \delta_k)$ .
- (iv) When the set  $\mathcal{T}_{k+1}$  contains many points, say more than some specified number  $N$ , one may consider an aggregation step similar to (iii) above and use  $AG(\mathcal{T}_{k+1}, \delta_k)$ .

To preserve convergence, we take  $\{\delta_k\}$  as a positive decreasing sequence tending to zero as  $k \rightarrow \infty$  (e.g.,  $\delta_k = \delta/(k+1)$  for some  $\delta > 0$ ) and, to avoid possible oscillations which may be caused by alternating local optimization in (ii) and aggregation of points in (iii), the local optimization is only used when there is an increase of  $\phi(\cdot)$  in two successive passages through step 1. The final algorithm is as follows. Remember that  $W_k^+$  denotes the vector formed by the strictly positive weights of  $\xi_k$  (in practice we only keep weights larger than some small constant, typically  $10^{-6}$ , all other points are removed from  $\xi_k$  and their total weight is reallocated to remaining points proportionally to the gradient components).

**Algorithm  $A_4$**  Let  $\{\gamma_k\}$  and  $\{\delta_k\}$ ,  $k \geq 0$ , denote two decreasing sequences of numbers in  $(0, 1)$ , with  $\lim_{k \rightarrow \infty} \gamma_k = \lim_{k \rightarrow \infty} \delta_k = 0$ .

**Step 0)** Choose  $X_0 = \{x^{(1)}, \dots, x^{(m)}\} \in \mathcal{X}^m$  such that  $M(\nu_0)$  has full rank, with  $\nu_0$  the uniform measure on  $X_0$ ; choose some arbitrary test-set  $\mathcal{T}_0$  of  $n$  points in  $\mathcal{X}$ ; choose  $N > n$  and  $\epsilon > 0$ , set  $k = 0$ ,  $\phi_{\text{old}} = \phi(\nu_0)$ .

**Step 1)** Compute  $\xi_k = A_0[X_k, \gamma_k \epsilon / 2]$ .

**Step 1a)** Compute  $S'_k = LM[\phi(\cdot | W_k^+); S_k]$ .

**Step 1b)** If  $\phi(S'_k | W_k^+) > \phi_{\text{old}}$ , substitute  $S'_k$  for  $S_k$  as support of  $\xi_k$ .

**Step 1c)** Compute  $\phi_k = \phi(\xi_k)$ , set  $\phi_{\text{old}} = \phi_k$ .

**Step 2)** Find  $x^* = \arg \max_{x \in \mathcal{T}_k} F(\xi_k, x)$  and compute  $x_k^* = LM[F(\xi_k, \cdot); x^*]$ .

**Step 3)** If  $F(\xi_k, x_k^*) \geq \epsilon/2$ , go to step 7; otherwise go to step 4.

**Step 4)** Find  $x^{**} = \arg \max_{x \in \mathcal{X}_\ell} F(\xi_k, x)$ .

**Step 5)** If  $F(\xi_k, x^{**}) < \epsilon/2$ , stop.

**Step 6)** Compute  $x_k^* = LM[F(\xi_k, \cdot); x^{**}]$ .

**Step 7)** Set  $X^* = EX[\{S_k, x_k^*\}; \delta_k]$  and  $\mathcal{T}^* = (\mathcal{T}_k, x_k^*)$ .

**Step 7a)** If  $\#\mathcal{T}^* > N$ , aggregate points in  $\mathcal{T}^*$ :  $\mathcal{T}_{k+1} = AG(\mathcal{T}^*, \delta_k)$ ; otherwise set  $\mathcal{T}_{k+1} = \mathcal{T}^*$ .

**Step 7b)** Aggregate points in  $X^*$ :  $X_{k+1} = AG(X^*, \delta_k)$ ,  $k \leftarrow k + 1$ , go to step 1.

In practice the finite set  $\mathcal{X}_\ell$  can be taken as an adaptive grid such that, as the number  $\ell$  of points tends to infinity, the minimax distance  $\rho = \max_{y \in \mathcal{X}} \min_{x \in \mathcal{X}_\ell} \|y - x\|_q$  tends to zero, for some  $L_q$ -norm  $\|\cdot\|_q$ . The search for  $x^{**} = \arg \max_{x \in \mathcal{X}_\ell} F(\xi_k, x)$  of step 4 is then replaced by a sequential inspection of grid points: the screening through the grid is stopped when either

- (i) a grid point  $x_i$  is found such that  $F(\xi_k, x) \geq \epsilon/2$ , this point being then taken as  $x^{**}$ ,
- or
- (ii) the value of  $\rho$  is small enough to ensure that  $\max_{x \in \mathcal{X}} F(\xi_k, x) < \epsilon$ , that is,  $\xi_k$  is  $\epsilon$ -optimal on  $\mathcal{X}$ .

This is detailed in Sect. 3.4.

### 3.2 Local maximizations

Steps 1a and 2 involve local maximizations, respectively with respect to the  $m_k$  points forming the support  $S_k$  of  $\xi_k$ ,  $S_k \subset \mathcal{X}^{m_k}$ , and with respect to  $x \in \mathcal{X}$ . There exist many sophisticated constrained optimization methods (such as sequential quadratic programming for instance) which can solve this problem efficiently. Note that very high accuracy is not required, any improvement compared with the initialization is enough to ensure convergence of algorithm A4. Simple unconstrained optimization methods can therefore also be used when  $\mathcal{X}$  has a simple form. Some methods are suggested below for situations where  $\mathcal{X}$  is the hypercube  $[-1, 1]^d$  and for the case of mixture experiments with  $\mathcal{X}$  given by the simplex  $\mathcal{P}_d$  defined by (5).

**Remark 3** *Note that a local maximization of  $F(\xi_k, x)$  initialized at the support points of  $\xi_k$  is generally not adequate at step 2. Consider for instance  $D$ -optimal design for the linear regression model  $\eta(x, \theta) = \theta_1 + \theta_2 x \sin(x)$  on  $\mathcal{X} = [0, \bar{x}]$ , with  $\bar{x} \geq 2$ . For  $\xi = (1/2)\delta_0 + (1/2)\delta_2$ ,  $F(\xi, 0) = F(\xi, 2) = 2$  and  $F(\xi, x)$  is locally maximum at  $x = 0$  and  $x = 2$  (with  $\xi$  being  $D$ -optimal on  $\mathcal{X}$  when  $2 \leq \bar{x} \leq \pi$ ), but  $\max_{x \in \mathcal{X}} F(\xi, x) > 2$ , and  $\xi$  is not  $D$ -optimal, when  $\bar{x} > \pi$ .*

**Reparameterization** When  $\mathcal{X} = [-1, 1]^d$ , one may consider the reparameterization  $x = \cos(y) \in \mathcal{X}$  (to be understood componentwise) with a local maximization with respect to  $y \in \mathbb{R}^d$  at step 2, or with respect to the  $m_k$  points  $y^{(1)}, \dots, y^{(m_k)}$ , each one in  $\mathbb{R}^d$ , at step 1a; see Appendix B for further details. Any classical method for unconstrained optimization can then be used, including derivative-free methods, see for instance Powell (1964).

**Projected gradient** Let  $h(\cdot)$  denote a differentiable function on some convex and compact set  $\mathcal{A}$ , with gradient  $\nabla h(\cdot)$ , and consider the local maximization of  $h(x)$  with respect to  $x \in \mathcal{A}$ . The projected-gradient algorithm uses iterations of the form

$$x_{j+1} = P_{\mathcal{A}} \left[ x_j^+(\alpha^*) \right]$$

with  $x_0$  given in  $\mathcal{A}$ , where  $x_j^+(\alpha) = x_j + \alpha \nabla h(x_j)$ ,

$$\alpha^* = \arg \max_{\alpha} h(P_{\mathcal{A}}[x_j^+(\alpha)])$$

and  $P_{\mathcal{A}}[\cdot]$  denotes the orthogonal projection onto  $\mathcal{A}$ .

When  $\mathcal{X} = [-1, 1]^d$ ,  $P_{\mathcal{X}}[x]$  simply amounts to a truncation of all components of  $x$  to  $[-1, 1]$ . The value  $\alpha^*$  can be obtained by a classical line-search algorithm, restricted to  $\alpha$  in the interval  $[0, \alpha_{\max}]$ , with

$$\alpha_{\max} = \min \left\{ \min_{i=1, \dots, d: \{\nabla h(x_j)\}_i \geq 0} (1 - \{x_j\}_i) / \{\nabla h(x_j)\}_i, \right. \\ \left. \min_{i=1, \dots, d: \{\nabla h(x_j)\}_i \leq 0} (1 + \{x_j\}_i) / |\{\nabla h(x_j)\}_i| \right\},$$

which guarantees that  $x_j^+(\alpha) \in \mathcal{X}$ . The same method can be applied when the optimization is with respect to  $S_k \in \mathcal{X}^{m_k}$ ; projections on  $\mathcal{X}^{m_k}$  then correspond to  $m_k$  projections over  $\mathcal{X}$ . The iterations can be stopped when the norm of the projected gradient is smaller than some specified threshold.

The projected-gradient method can also be used in the case of mixture experiments where  $\mathcal{X}$  is the simplex  $\mathcal{P}_d$  given by (5). The projection of  $x_j^+(\alpha)$  on  $\mathcal{X}$  can then be obtained as the solution of a Quadratic Programming (QP) problem:  $P_{\mathcal{P}_d}[x_j^+(\alpha)]$  minimizes  $\|x - x_j^+(\alpha)\|^2$  with respect to  $x$  satisfying the linear constraints  $x \in \mathcal{P}_d$ . Alternatively, one may also use the property that, for any  $Z \in \mathbb{R}^d$ , the projection  $P_{\mathcal{P}_d}(Z)$  is given by  $x(Z, t^*)$ , where  $x(Z, t) = \max\{Z - t1_d, 0_d\}$  (componentwise, with  $1_d$  and  $0_d$  respectively the  $d$ -dimensional vectors of ones and zeros) and  $t^*$  maximizes  $L[x(Z, t), t]$  with respect to  $t$ , with  $L(x, t)$  the partial Lagrangian  $L(x, t) = (1/2)\|x - Z\|^2 + t[1_d^\top x - 1]$ . Indeed,  $L(x, t)$  can be written as

$$L(x, t) = (1/2)\|x - (Z - t1_d)\|^2 + t[1_d^\top Z - 1] - dt^2/2,$$

which reaches its minimum with respect to  $x \in \mathbb{R}^{d+}$  for  $x = x(Z, t)$ . One may notice that  $\max\{\max_i\{Z\}_i - t^*, 0\} \leq 1_d^\top x(Z, t^*) = 1 \leq d \max\{\max_i\{Z\}_i - t^*, 0\}$ , so that the search for  $t^*$  can be restricted to the interval  $[\max_i\{Z\}_i - 1, \max_i\{Z\}_i - 1/d]$ , see Pronzato and Pázman (2013, Chap. 9).



### 3.3 From $(\epsilon/2)$ -optimality on a grid to $\epsilon$ -optimality on a compact set

For any  $\xi \in \Xi(\mathcal{X})$  and any  $x, y \in \mathcal{X}$  we have

$$|F(\xi, x) - F(\xi, y)| = |\text{trace}\{\mu(x) - \mu(y)\nabla\Phi[M(\xi)]\}| \leq \|\mu(x) - \mu(y)\| \|\nabla\Phi[M(\xi)]\|.$$

Only the second term on the right-hand side depends on the criterion  $\phi(\cdot)$  and design  $\xi$ , the first term being related to the problem under consideration. For instance (assuming that  $M(\xi)$  has full rank),  $D$ -optimality with  $\phi(\xi) = \log \det[M(\xi)]$  gives  $\nabla\Phi[M(\xi)] = M^{-1}(\xi)$ ,  $A$ -optimality with  $\phi(\xi) = -\text{trace}[M^{-1}(\xi)]$  gives  $\nabla\Phi[M(\xi)] = M^{-2}(\xi)$  and, more generally, when  $t \geq 1$  and  $A$  is non-negative definite  $\phi_{A,t}(\xi) = -\text{trace}[AM^{-t}(\xi)]$  gives  $\nabla\Phi[M(\xi)] = \sum_{i=1}^t M^{-i}(\xi)A^\top M^{-(t+1-i)}(\xi)$ , see (Harville, 1997, Chap. 15).

Suppose that  $\mu(\cdot)$  satisfies the Lipschitz condition

$$\forall y \in \mathcal{X} \cap \mathcal{B}_q(x, \rho), \quad \|\mu(y) - \mu(x)\| \leq L_\mu(x, \rho) \rho,$$

with  $\mathcal{B}_q(x, \rho)$  the ball  $\{y : \|y - x\|_q \leq \rho\}$  for some norm  $\|\cdot\|_q$ ,  $1 \leq q \leq \infty$ . We then obtain

$$\forall y \in \mathcal{X} \cap \mathcal{B}_q(x, \rho), \quad F(\xi, y) \leq F(\xi, x) + \rho L_F(\xi, x, \rho), \quad (9)$$

where  $L_F(\xi, x, \rho) = L_\mu(x, \rho) \|\nabla\Phi[M(\xi)]\|$ . Suppose in particular that  $\mu(\cdot)$  is differentiable on  $\mathcal{X}$  which is a convex set. Then, for any  $x, y$  in  $\mathcal{X}$  we can write

$$F(\xi, y) = F(\xi, x) + \left. \frac{\partial F(\xi, z)}{\partial z^\top} \right|_{z=x+\alpha(y-x)} (y - x)$$

for some  $\alpha \in [0, 1]$ , so that (9) is satisfied for  $q = 2$  with

$$L_F(\xi, x, \rho) = \max_{\max z \in \mathcal{B}_2(x, \rho)} \left\| \frac{\partial F(\xi, z)}{\partial z} \right\|.$$

More generally, we can use regularity properties of  $\mu(\cdot)$  to derive inequalities of the form

$$\forall y \in \mathcal{X} \cap \mathcal{B}_q(x, \rho), \quad F(\xi, y) \leq F(\xi, x) + h(\xi, x, \rho), \quad (10)$$

where, for any  $\xi$  such that  $M(\xi)$  has full rank,  $\sup_{x \in \mathcal{X}} h(\xi, x, \rho) \rightarrow 0$  as  $\rho \rightarrow 0$ .

Consider then the final design  $\xi_k$  returned by algorithm  $A_4$  and take  $\rho = \max_{y \in \mathcal{X}} \min_{x \in \mathcal{X}_\ell} \|y - x\|_q$  with  $\mathcal{X}_\ell$  the set used at step 4 of the algorithm. Then (10) implies that

$$\max_{y \in \mathcal{X}} F(\xi, y) \leq \max_{x \in \mathcal{X}_\ell} \{F(\xi_k, x) + h(\xi, x, \rho)\}$$

and, if  $\rho$  is small enough to ensure that  $\max_{x \in \mathcal{X}_\ell} h(\xi_k, x, \rho) < \epsilon/2$ ,  $\xi_k$  is  $\epsilon$ -optimal on  $\mathcal{X}$  and the efficiency bound (7) applies.

The case where  $\mu(x)$  has rank one, *i.e.*,  $\mu(x) = f(x)f^\top(x)$  for some  $p$ -dimensional regressor vector  $f(x)$ , deserves special attention. Consider  $\phi_0(\xi) = \log \det[M(\xi)]$  ( $D$ -optimality) and  $\phi_t(\xi) = -\text{trace}[M^{-t}(\xi)]$  for  $t \geq 1$ . For each of these criteria we get

$$\begin{aligned} |F(\xi, y) - F(\xi, x)| &= c_t \left| f^\top(y)M^{-(t+1)}(\xi)f(y) - f^\top(x)M^{-(t+1)}(\xi)f(x) \right| \\ &= c_t \left| \Delta f^\top M^{-(t+1)}(\xi)\Delta f + 2\Delta f^\top M^{-(t+1)}(\xi)f(x) \right| \\ &\leq c_t \frac{\|\Delta f\|^2}{\lambda_{\min}^{t+1}[M(\xi)]} + 2c_t \|f(x)\| \frac{\|\Delta f\|}{\lambda_{\min}^{t+1}[M(\xi)]}, \end{aligned} \quad (11)$$

where we have denoted  $c_t = \max(t, 1)$  and  $\Delta f = f(y) - f(x)$ . This gives (10) with

$$h(\xi, x, \rho) = \frac{c_t}{\lambda_{\min}^{t+1}[M(\xi)]} [D_\rho^2(x) + 2D_\rho(x)\|f(x)\|] \quad (12)$$

where  $D_\rho(x) = \max_{y \in \mathcal{X} \cap \mathcal{B}_q(x, \rho)} \|f(y) - f(x)\|$ .

If  $f(\cdot)$  is differentiable, then we also have  $\partial F(\xi, x)/\partial x = 2c_t f^\top(x)M^{-(t+1)}(\xi)\partial f(x)/\partial x$  so that

$$\left\| \frac{\partial F(\xi, x)}{\partial x} \right\| \leq 2c_t \|f(x)\| \left\| \frac{\partial f(x)}{\partial x} \right\| \frac{1}{\lambda_{\min}^{t+1}[M(\xi)]}.$$

This gives (9) with  $L_F(\xi, x, \rho) = \{2c_t/\lambda_{\min}^{t+1}[M(\xi)]\} \max_{y \in \mathcal{X} \cap \mathcal{B}_2(x, \rho)} \{\|f(y)\| \|\partial f(y)/\partial y\|\}$ . Any particular problem can thus be handled by a case-by-case analysis in order to obtain a Lipschitz inequality similar to (9) or (10). Examples with polynomial regression models are considered in Appendix A.

### 3.4 Adaptive grids

An adaptive construction of the grid  $\mathcal{X}_\ell$  used at step 4 of algorithm  $A_4$  allows us to only refine the grid at locations where the upper bound (10) needs to be improved. Meanwhile, it also sometimes speeds up the search for  $x^{**}$ . Here we indicate a possible construction when  $\mathcal{X} = [-1, 1]^d$ .

For any  $n \geq 0$ , let  $X_{1,mm}(n) = \{x^{(1)}, \dots, x^{(n)}\}$  denote the design (sometimes called minimax-distance optimal) given by  $x^{(j)} = (2j-1)/n-1$ ,  $j = 1, \dots, n$ , and  $X_{d,mm}(n) = X_{1,mm}^{\otimes d}(n)$  denote the  $n^d$ -point  $d$ -dimensional grid with all coordinates in  $X_{1,mm}(n)$ .

Step 4 of algorithm  $A_4$  is then replaced by the following (step 5 is removed).

**Step 4-0)** Choose  $n$  (say,  $n = 100$ ),  $N_{\max}$  (say,  $N_{\max} = 10^7$ ), set  $\mathcal{X}_\ell^{(0)} = X_{d,mm}(n)$ ,  $G_0 = n^d$ ,  $\rho_0 = 1/m$ ,  $j = 0$ .

**Step 4-1)** Find  $x^{**} = \arg \max_{x \in \mathcal{X}_\ell^{(j)}} F(\xi_k, x)$ .

**Step 4-2)** If  $F(\xi_k, x^{**}) \geq \epsilon/2$ , go to step 6 of  $A_4$ .

$p$	support points					
2	-1	1				
3	-1	0	1			
4	-1	$-\frac{1}{\sqrt{5}}$	$\frac{1}{\sqrt{5}} \simeq 0.4472$	1		
5	-1	$-\frac{\sqrt{3}}{\sqrt{7}}$	0	$\frac{\sqrt{3}}{\sqrt{7}} \simeq 0.6547$	1	
6	-1	$-\frac{\sqrt{7+2\sqrt{7}}}{\sqrt{21}}$	$-\frac{\sqrt{7-2\sqrt{7}}}{\sqrt{21}}$	$\frac{\sqrt{7-2\sqrt{7}}}{\sqrt{21}} \simeq 0.2852$	$\frac{\sqrt{7+2\sqrt{7}}}{\sqrt{21}} \simeq 0.7651$	1

Table 1:  $D$ -optimal designs for polynomial regression on  $[-1, 1]$  (all support points are equally weighted).

**Step 4-3)** If  $\bar{\epsilon} = \max_{x \in \mathcal{X}_\ell^{(j)}} F(\xi_k, x) + h(\xi_k, x, \rho_j) < \epsilon$  or if  $G_j > N_{\max}$ , stop  $A_4$ :  $\xi_k$  is  $\bar{\epsilon}$ -optimal on  $\mathcal{X}$ .

**Step 4-4)** Set  $\mathcal{X}_\ell^{(j>)} = \{x^{(j)} \in \mathcal{X}_\ell^{(j)} : F(\xi_k, x^{(j)}) + h(\xi_k, x^{(j)}, \rho_j) \geq \epsilon\}$ ; for each  $x^{(j)} \in \mathcal{X}_\ell^{(j>)}$ , divide the hypercube  $\mathcal{B}_\infty(x^{(j)}, \rho_j)$  into  $2^d$  hypercubes  $\mathcal{B}_\infty(x^{(ji)}, \rho_j/2)$ , where  $x^{(ji)} = x^{(j)} + \rho_j x_{d,mm}^{(i)}(2)$  with  $x_{d,mm}^{(i)}(2)$  the  $i$ -th element of the  $2^d$ -point grid  $X_{d,mm}(2)$ ,  $i = 1, \dots, 2^d$ . Collect all  $x^{(ji)}$  to form  $\mathcal{X}_\ell^{(j+1)}$ , set  $G_{j+1} = \#\mathcal{X}_\ell^{(j+1)}$ ,  $\rho_{j+1} = \rho_j/2$ ,  $j \leftarrow j+1$ , go to step 4-1.

Steps 4-3 and 4-4 rely on (10); although the construction of the grids  $\mathcal{X}_\ell^{(j)}$  use the norm  $\|\cdot\|_\infty$ , another norm can be used in the definition of  $h(\xi, x, \rho)$ . Due to the bound  $N_{\max}$  set on the number  $G_j$  of elements of  $\mathcal{X}_\ell^{(j)}$ , the algorithm may stop before the precision on  $\max_{x \in \mathcal{X}} F(\xi_k, x)$  reaches  $\epsilon$ , *i.e.*, one may have  $\bar{\epsilon} > \epsilon$  at step 4-3. The use of a precise bound in (10) is then crucial to avoid a fast increase of  $G_j$  and a premature stopping of the algorithm.

## 4 Examples

### 4.1 Description

We have tested algorithm  $A_4$  on a series of examples with known optimal design  $\xi^*$ . The design space is  $\mathcal{X} = [-1, 1]^d$  with  $d = 1$  or  $2$  and the number of parameters  $p$  varies between 3 and 9. We use  $\phi(\xi) = \log \det M(\xi)$  for  $D$ -optimal design and  $\phi(\xi) = -\text{trace}[M^{-1}(\xi)]$  for  $A$ -optimal design.

Problems 1 to 4 correspond to  $D$ -optimal design for linear regression models with  $f(x)$  given by (13) and  $p = s$  varying from 3 to 6. For each  $p$ , the  $D$ -optimal design measure is uniquely defined and gives weight  $1/p$  at the roots of the polynomial  $t \mapsto (1 - t^2)P'_{p-1}(t)$  with  $P'_k(\cdot)$  the derivative of the  $k$ -th Legendre polynomial, see Table 1.

Problems 5 and 6 correspond to  $D$ -optimal design for additive polynomial models:  $f(x)$  is given by (14) with  $d = 2$ ,  $p = 2s - 1$  and  $s = 3, 4$  respectively. Problem 7 correspond to  $D$ -optimal design for a complete product-type interaction model with  $d = 2$ :  $f(x) = g(x_1) \otimes g(x_2)$

with  $\otimes$  denoting tensor product and  $g(x)$  given by (13) with  $s = 3$ ; this gives  $p = 9$ . For those three problems the  $D$ -optimal design is the cross-product of the  $D$ -optimal design measures for the individual models; see Schwabe (1996, Chap. 4 and 5).

In Problems 8 and 9, the models are the same as in Problems 1 and 7 but the criterion is  $\phi(\xi) = -\text{trace}[M^{-1}(\xi)]$  ( $A$ -optimality). The  $A$ -optimal design for Problem 8 is

$$\xi_A^* = \begin{Bmatrix} -1 & 0 & 1 \\ 1/4 & 1/2 & 1/4 \end{Bmatrix},$$

the  $A$ -optimal design for Problem 9 is the product measure  $\xi_A^* \otimes \xi_A^*$ ; see Schwabe (1996, Chap. 4).

Problem 10 corresponds to (local)  $D$ -optimal design for the nonlinear regression model

$$\eta(x, \theta) = \theta_0 + \theta_1 \exp(-\theta_2 x_1) + \frac{\theta_3}{\theta_3 - \theta_4} [\exp(-\theta_4 x_2) - \exp(-\theta_3 x_2)]$$

with five parameters  $\theta = (\theta_0, \theta_1, \theta_2, \theta_3, \theta_4)$  and two design variables  $x = (x_1, x_2) \in \mathcal{X} = [0, 2] \times [0, 10]$ . A numerical value must be given to the parameters  $\theta_2$ ,  $\theta_3$  and  $\theta_4$  that intervene nonlinearly in  $\eta(x, \theta)$ , and we use  $\theta_2 = 2$ ,  $\theta_3 = 0.7$ ,  $\theta_4 = 0.2$ . The model considered is additive with a constant term, so that the tensor product of the  $D$ -optimal designs for the two models  $\eta_1(x_1, \beta^{(1)}) = \beta_0^{(1)} + \beta_1^{(1)} \exp(-\beta_2^{(1)} x_1)$  and  $\eta_2(x_2, \beta^{(2)}) = \beta_0^{(2)} + \beta_1^{(2)} [\exp(-\beta_2^{(2)} x_2) - \exp(-\beta_1^{(2)} x_2)] / (\beta_1^{(2)} - \beta_2^{(2)})$ , respectively at  $\beta^{(1)} = (\theta_0, \theta_1, \theta_2)$  and  $\beta^{(2)} = (\theta_0, \theta_3, \theta_4)$ , is  $D$ -optimal for  $\eta(x, \theta)$  at  $\theta$ , see Schwabe (1995). These two  $D$ -optimal designs are supported on three points only, each one receiving mass  $1/3$ , which can be computed with arbitrary precision using Theorem 1-(iii): the support points are approximately 0, 0.46268527927 and 2 for  $\eta_1(x_1, \beta^{(1)})$  and 0, 1.22947139883, 6.85768905493 for  $\eta_2(x_2, \beta^{(2)})$ . The design space  $\mathcal{X}$  is renormalized to  $[-1, 1]^2$  in order to use the adaptive-grid construction of Sect. 3.4; the tests for optimality are based on (10) and (12).

## 4.2 Results and discussion

The results obtained are presented in Table 2. When  $d = 1$ , the algorithm is initialized with  $X_0$  having  $m$  equidistant points,  $X_0 = \{(2i - 1)/m - 1, i = 1 \dots, m\}$ , with  $m = p(p + 1)/2$ ; when  $d = 2$ ,  $X_0$  is the tensor product of two such designs, with  $m^2$  the smallest integer larger than or equal to  $p(p + 1)/2$ . The test-set  $\mathcal{T}_0$  is the union of  $X_0$  and the  $2^d$  vertices of  $\mathcal{X}$ . The number of points in  $\mathcal{T}_k$  always remains reasonably small (see Table 2) so that step 7a is not used. The decreasing sequences  $\{\gamma_k\}$  and  $\{\delta_k\}$  are respectively given by  $1/(k + 1)$  and  $0.1/(k + 1)$ ,  $\epsilon$  is set to  $10^{-6}$  in all examples. The algorithm  $A_0$  used at step 1 is based on a combination of projected-gradient and vertex-exchange methods, see Wu (1978a). Local maximizations at steps 1a and 2 of  $A_4$  use reparameterization and the derivative-free algorithm of Powell (1964).

The set  $\mathcal{X}_\ell$  of step 4 corresponds to an adaptive grid, constructed according to the algorithm of Sect. 3.4 with  $n = 100$  for  $d = 1$  and  $n = 10$  for  $d = 2$  (so that  $G_0 = 100$ ).

In all problems considered the optimal design  $\xi^*$  is unique and one can thus compute its distance to  $\xi_{k_{\max}}$  returned by  $A_4$ . Various metrics can be considered and we use here the Wasserstein and Lévy-Prokhorov metrics, see Appendix C. Table 2 indicates that  $\xi_{k_{\max}}$  is very close to  $\xi^*$  for all problems considered. We also report in Table 2 the distance to 1 of the efficiencies  $\mathcal{E}_\phi(\xi_{k_{\max}})$  defined by (2) together with the values of  $1 - \underline{\mathcal{E}}_\phi(\xi_{k_{\max}})$ , with  $\underline{\mathcal{E}}_\phi(\xi_{k_{\max}})$  the efficiency bound given by (7). Note that  $1 - \underline{\mathcal{E}}_\phi(\xi) < \epsilon/p$  for  $D$ -optimality when  $\max_{x \in \mathcal{X}} F(\xi, x) < \epsilon$  and  $\phi(\xi) = \log \det M(\xi)$  and  $1 - \underline{\mathcal{E}}_\phi(\xi) < \epsilon/\text{trace}[M^{-1}(\xi)]$  for  $A$ -optimality when  $\max_{x \in \mathcal{X}} F(\xi, x) < \epsilon$  and  $\phi(\xi) = -\text{trace}[M^{-1}(\xi)]$ . Due to the fact that the efficiency bound, used to terminate the algorithm, is very pessimistic in most problems, the true efficiency of  $\xi_{k_{\max}}$  is often much closer to 1 than the bound indicated. We thus observe in all cases considered that a very good precision is achieved although the number of iterations  $k_{\max}$  is small. Also note that the construction of the adaptive grid is only required a small number of times (only once for more than half of the problems considered); due to the local maximizations performed at steps 1a, 2 and 6, the points in the test-set  $\mathcal{T}_k$ , although there are only few of them, are generally enough to construct an optimal design on the whole set  $\mathcal{X}$ .

Most of the computational cost of  $A_4$  corresponds to step 1 and especially step 4.

Concerning step 4 (which we try to avoid as much as possible), the use of an adaptive grid, as proposed in Sect. 3.4, aims at minimizing its cost. Using precise bounds in the construction of  $h(\xi, x, \delta)$  used in (10) is essential to maintain the cardinality  $G_j$  of  $\mathcal{X}_\ell^{(j)}$  reasonably small. One may notice that the construction of Sect. 3.4 does not make use of all the information contained in the evaluations of  $F(\xi_k, x_i)$  for the various  $x_i$  which are used. Indeed, some  $x^{(j)}$  are removed from  $\mathcal{X}_\ell^{(j)}$  when constructing the set  $\mathcal{X}_\ell^{(j>)}$ . With some computational effort, other constructions making a more efficient use of the information collected on  $F(\xi_k, \cdot)$  could be considered, and the results in (Harman and Pronzato, 2007; Pronzato, 2013) could be used to remove parts of  $\mathcal{X}$  with low values of  $F(\xi_k, x)$  that cannot contain support points of an optimal design. Also, Lipschitz bounds obtained from diagonal partitions, see Sergeyev and Kvasov (2006), might be used to get a significant reduction of the number of evaluations of  $F(\xi_k, x_i)$  required when  $d > 1$ . Of course, for a given problem this number depends very much on the value of  $\epsilon$ , see the bottom part of Table 2.

Concerning step 1, the algorithm is constructed in such a way that  $X_k$  has a small number of elements. The choice of the optimization method used for  $A_0$  at step 1 is then not crucial. We have used a combination of the projected-gradient and vertex-exchange methods for the results in Table 2. Numerical experimentations with a pure projected-gradient algorithm (McCormick and Tapia, 1972), the methods of cutting planes (Kelley, 1960) and the level method (Nesterov,

2004, Sect. 3.3.3), originally developed for non-differentiable optimization (see also Pronzato and Pázman (2013, Chap. 9)), yield similar results.

## 5 Conclusion

We have proposed an algorithm for the construction of optimal design measures over a compact set  $\mathcal{X}$  for a concave and differentiable criterion  $\phi(\cdot)$ . The method exploits the property that optimal designs are often supported on a small number of points and combines updating of the support with convex optimization of the weights it receives. Efficiency bounds and guaranteed  $\epsilon$ -optimality over  $\mathcal{X}$  are obtained using Lipschitz continuity of the directional derivative of  $\phi(\cdot)$ . As illustrated by some examples, the algorithm seems to be quite effective when accurate Lipschitz constants can be determined for the particular problem considered. General techniques for deriving those constants are given which could be applied to any model satisfying usual regularity conditions. Also, some general ideas have been given for the local optimization of the support of a design measure, which should allow the construction of specific methods especially adapted to particular design regions  $\mathcal{X}$ .

On the other hand, some issues remain that call for further developments. We mention two of them.

(i) There exist problems for which the optimal design is not unique; in such situations it would be of interest to have an algorithm capable of ensuring convergence to an  $\epsilon$ -optimal design with minimal support. This is not the case of our algorithm  $A_4$  and improvements in that direction are under current investigation. Note that once an optimal design is found, the determination of another optimal design (having necessarily the same information matrix when  $\Phi(\cdot)$  is strictly concave) with minimum support corresponds to a fixed-charge problem, see, *e.g.*, Sadagopan and Ravindran (1982). A simple reduction of the support of an optimal design can be obtained by iteratively removing elementary matrices  $\mu(x^{(i)})$  that are in the convex hull of others  $\mu(x^{(j)})$ ,  $j \neq i$ .

(ii) Although the convex optimization of weights for non-differentiable criteria (for instance,  $E$ -optimality) is not a big challenge, see, *e.g.*, Nesterov (2004, Sect. 3.3.3), Pronzato and Pázman (2013, Chap. 9), the determination of new points being candidates for inclusion in the support and the construction of efficiency bounds over  $\mathcal{X}$  raise specific issues due to the fact that the maximum of the directional derivative of  $\phi(\cdot)$  is not always attained at a one-point (delta) measure. The geometrical characterization of optimal designs, see in particular Dette and Studden (1993), may then prove salutary for the construction of efficient algorithms.

Pb.	$\epsilon$	$d$	$p$	$k_{\max}$	# step 4	$W_1$	$W_2$	$L$	$1 - \mathcal{E}_\phi(\xi_{k_{\max}})$	$1 - \underline{\mathcal{E}}_\phi$	$\#\mathcal{T}_{k_{\max}}$	$G_{\max}$
1	$10^{-6}$	1	3	4	1	$4.4 \cdot 10^{-10}$	$7.6 \cdot 10^{-10}$	$1.3 \cdot 10^{-9}$	$1.3 \cdot 10^{-11}$	$2.8 \cdot 10^{-7}$	8	195
2	$10^{-6}$	1	4	5	4	$3.4 \cdot 10^{-5}$	$5.7 \cdot 10^{-5}$	$1.1 \cdot 10^{-4}$	$1.6 \cdot 10^{-8}$	$1.5 \cdot 10^{-7}$	13	402
3	$10^{-6}$	1	5	6	1	$2.4 \cdot 10^{-6}$	$4.7 \cdot 10^{-6}$	$1.1 \cdot 10^{-5}$	$2.5 \cdot 10^{-10}$	$1.3 \cdot 10^{-7}$	15	672
4	$10^{-6}$	1	6	5	1	$9.1 \cdot 10^{-6}$	$1.4 \cdot 10^{-5}$	$2.6 \cdot 10^{-5}$	$2.4 \cdot 10^{-9}$	$1.0 \cdot 10^{-7}$	12	924
5	$10^{-6}$	2	5	11	2	$4.5 \cdot 10^{-8}$	$1.8 \cdot 10^{-4}$	$7.3 \cdot 10^{-8}$	$4.2 \cdot 10^{-15}$	$1.5 \cdot 10^{-7}$	21	$\simeq 2.1 \cdot 10^3$
6	$10^{-6}$	2	7	20	2	$4.2 \cdot 10^{-8}$	$9.7 \cdot 10^{-5}$	$7.1 \cdot 10^{-8}$	$1.2 \cdot 10^{-15}$	$8.7 \cdot 10^{-8}$	42	$\simeq 1.4 \cdot 10^6$
7	$10^{-6}$	2	9	7	3	$5.7 \cdot 10^{-9}$	$2.3 \cdot 10^{-8}$	$3.5 \cdot 10^{-8}$	$1.8 \cdot 10^{-11}$	$8.8 \cdot 10^{-8}$	55	$\simeq 1.1 \cdot 10^6$
8	$10^{-6}$	1	3	10	1	$9.9 \cdot 10^{-9}$	$8.9 \cdot 10^{-5}$	$5.2 \cdot 10^{-9}$	$< 10^{-16}$	$8.4 \cdot 10^{-8}$	8	196
9	$10^{-6}$	2	9	13	1	$1.0 \cdot 10^{-8}$	$5.7 \cdot 10^{-5}$	$2.2 \cdot 10^{-8}$	$1.3 \cdot 10^{-15}$	$1.2 \cdot 10^{-8}$	53	$\simeq 4.9 \cdot 10^6$
10	$10^{-6}$	2	5	9	1	$2.8 \cdot 10^{-7}$	$2.1 \cdot 10^{-4}$	$3.6 \cdot 10^{-8}$	$4.8 \cdot 10^{-14}$	$1.8 \cdot 10^{-7}$	21	$\simeq 6.4 \cdot 10^6$
10	$10^{-5}$	2	5	7	1	$1.1 \cdot 10^{-6}$	$7.7 \cdot 10^{-4}$	$1.5 \cdot 10^{-6}$	$1.9 \cdot 10^{-12}$	$1.9 \cdot 10^{-6}$	21	$\simeq 1.7 \cdot 10^6$
10	$10^{-4}$	2	5	7	1	$5.0 \cdot 10^{-7}$	$1.9 \cdot 10^{-6}$	$1.6 \cdot 10^{-6}$	$1.6 \cdot 10^{-12}$	$1.3 \cdot 10^{-5}$	21	$\simeq 5.9 \cdot 10^5$
10	$10^{-3}$	2	5	6	1	$4.9 \cdot 10^{-6}$	$7.6 \cdot 10^{-6}$	$1.6 \cdot 10^{-5}$	$1.7 \cdot 10^{-10}$	$1.0 \cdot 10^{-4}$	21	$\simeq 2.2 \cdot 10^5$
10	$10^{-2}$	2	5	4	1	$1.6 \cdot 10^{-4}$	$2.4 \cdot 10^{-4}$	$6.1 \cdot 10^{-4}$	$1.8 \cdot 10^{-7}$	$1.5 \cdot 10^{-3}$	21	$\simeq 8.7 \cdot 10^4$

Table 2: Results obtained by using algorithm  $A_4$  on Problems 1-10 with  $\epsilon = 10^{-6}$  and on Problem 10 with  $\epsilon$  between  $10^{-2}$  and  $10^{-6}$ : number  $k_{\max}$  of iterations before stopping and number of passages through step 4; Wasserstein distances  $W_1(\xi_{k_{\max}}, \xi^*)$  and  $W_2(\xi_{k_{\max}}, \xi^*)$  and Lévy-Prokhorov distance  $L(\xi_{k_{\max}}, \xi^*)$  between  $\xi_{k_{\max}}$  and the optimal design  $\xi^*$ ; (1-)  $\phi$ -efficiency  $\mathcal{E}_\phi(\xi_{k_{\max}})$  and (1-) efficiency bound  $\underline{\mathcal{E}}_\phi(\xi_{k_{\max}})$  given by (7); number of elements of  $\mathcal{T}_{k_{\max}}$  and maximum number of points in the grids  $\mathcal{X}_\ell^{(j)}$  used to guarantee the  $\epsilon$ -optimality of  $\xi_{k_{\max}}$ .

## Appendix A: Bounds for directional derivatives in polynomial regression

In this appendix, we construct bound for  $\|f(x)\|$  and  $\|f(y) - f(x)\|$  to be used in (12) for polynomial regression models. Here  $\mathcal{X} = [-1, 1]$  and  $\mu(x) = f(x)f^\top(x)$  with

$$f(x) = (1, x, x^2, \dots, x^{s-1})^\top, \quad s \geq 2. \quad (13)$$

We then have the global Lipschitz bound  $\|\mu(x) - \mu(y)\| \leq \|A_s\| |x - y|$ , with  $A_s$  the  $s \times s$  matrix

$$A_s = \begin{pmatrix} 0 & 1 & 2 & \cdots & s-1 \\ 1 & 2 & & & \\ 2 & & \ddots & & \vdots \\ \vdots & & & & \\ s-1 & \cdots & & & 2(s-1) \end{pmatrix}$$

so that  $\|A_s\|^2 = \text{trace}(A_s^\top A_s) = s^2(7s-5)(s-1)/6$ . Also, for all  $x, y \in \mathcal{X}$ ,  $\|\Delta f\|^2 = \|f(y) - f(x)\|^2 \leq (x-y)^2 [1+2^2+\cdots+(s-1)^2] = (x-y)^2 s(s-1)(2s-1)/6$ , and  $\max_{x \in \mathcal{X}} \|f(x)\| = \sqrt{s}$ , which can be used as global bounds in (11). More precise local bounds can also be obtained:

$$\begin{aligned} \|\Delta f\|^2 &= \sum_{j=0}^{s-1} (y^j - x^j)^2 = (x-y)^2 \sum_{j=1}^{s-1} \left( \sum_{a+b=j-1} x^a y^b \right)^2 \\ &\leq (x-y)^2 \sum_{j=1}^{s-1} \left( \sum_{a+b=j-1} [\max\{|x|, |y|\}]^{a+b} \right)^2 = (x-y)^2 \omega_s(\max\{|x|, |y|\}) \end{aligned}$$

where  $\omega_s(c) = \sum_{j=1}^{s-1} j^2 c^{2(j-1)}$ . Therefore,

$$D_\delta(x) = \max_{y \in \mathcal{X} \cap \mathcal{B}_1(x, \delta)} \|f(y) - f(x)\| \leq \delta \sqrt{\omega_s(\max\{|x+\delta|, |x-\delta|\})}.$$

to be used in (10, 12).

Multi-factor polynomial regression models can be handled similarly. For additive models with

$$f(x) = (1, x_1, x_1^2, \dots, x_1^{s-1}, x_2, x_2^2, \dots, x_d, x_d^2, \dots, x_d^{s-1})^\top, \quad (14)$$

$x = (x_1, \dots, x_d) \in \mathcal{X} = [-1, 1]^d$ , then  $\|\Delta f\|^2 \leq \|x-y\|^2 s(s-1)(2s-1)/6$  and  $\max_{x \in \mathcal{X}} \|f(x)\| = [1 + d(s-1)]^{1/2}$ . We also obtain, similarly to the case  $d = 1$  considered above,

$$\|\Delta f\|^2 = \sum_{i=1}^d (y_i - x_i)^2 \sum_{j=1}^{s-1} \left( \sum_{a+b=j-1} x_i^a y_i^b \right)^2 \leq \sum_{i=1}^d (x_i - y_i)^2 \omega_s(\max\{|x_i|, |y_i|\})$$



so that

$$D_\delta(x) = \max_{y \in \mathcal{X} \cap \mathcal{B}_\infty(x, \delta)} \|f(y) - f(x)\| \leq \delta \left[ \sum_{i=1}^d \omega_s(\max\{|x_i + \delta|, |x_i - \delta|\}) \right]^{1/2}.$$

Consider now models with complete product-type interactions and take  $f(x) = g(x_1) \otimes g(x_2)$  with  $g(x)$  given by (13),  $x = (x_1, x_2) \in \mathcal{X} = [-1, 1]^2$ . Then  $\Delta f = f(y) - f(x) = [g(x_1) \otimes g(x_2)] - [g(y_1) \otimes g(y_2)] = [g(x_1) - g(y_1)] \otimes g(x_2) + g(y_1) \otimes [g(x_2) - g(y_2)]$  and

$$\begin{aligned} \|\Delta f\| &\leq \| [g(x_1) - g(y_1)] \otimes g(x_2) \| + \| g(y_1) \otimes [g(x_2) - g(y_2)] \| \\ &= \| [g(x_1) - g(y_1)] g^\top(x_2) \| + \| g(y_1) [g(x_2) - g(y_2)]^\top \| \\ &= \|g(x_1) - g(y_1)\| \|g(x_2)\| + \|g(y_1)\| \|g(x_2) - g(y_2)\| \\ &\leq \left( \frac{s(s-1)(2s-1)}{6} \right)^{1/2} \sqrt{s} [|x_1 - y_1| + |x_2 - y_2|] \\ &\leq s \frac{\sqrt{(s-1)(2s-1)}}{\sqrt{3}} \|x - y\|, \\ \|f(x)\| &= \|g(x_1)\| \|g(x_2)\| \leq s. \end{aligned}$$

We also get the local bounds

$$\begin{aligned} D_\delta(x) &= \max_{y \in \mathcal{X} \cap \mathcal{B}_\infty(x, \delta)} \|f(y) - f(x)\| \\ &\leq \delta \min \left\{ \|g(x_2)\| \sqrt{\omega_s(\max\{|x_1 + \delta|, |x_1 - \delta|\})} + \sqrt{s} \sqrt{\omega_s(\max\{|x_2 + \delta|, |x_2 - \delta|\})}, \right. \\ &\quad \left. \sqrt{s} \sqrt{\omega_s(\max\{|x_1 + \delta|, |x_1 - \delta|\})} + \|g(x_1)\| \sqrt{\omega_s(\max\{|x_2 + \delta|, |x_2 - \delta|\})} \right\}. \end{aligned}$$

## Appendix B: Reparameterization for optimization on an hypercube

Let  $h(\cdot)$  denote a function to be maximized with respect to  $x \in \mathcal{X} = [-1, 1]^d$ ,  $d \geq 1$ . We suppose that  $h(\cdot)$  is twice continuously differentiable in  $\mathcal{X}$ . Consider the reparameterization defined by  $x = x(y) = \cos(y)$  (componentwise), with  $y \in \mathbb{R}^d$ , and denote  $\tilde{h}(y) = h[x(y)]$ , so that

$$\frac{\partial \tilde{h}(y)}{\partial y} = -D(y) \frac{\partial h(x)}{\partial x} \Big|_{x=x(y)} \quad \text{and} \quad \frac{\partial^2 \tilde{h}(y)}{\partial y \partial y^\top} = D(y) \frac{\partial^2 h(x)}{\partial x \partial x^\top} \Big|_{x=x(y)} D(y),$$

where  $D(y)$  is the diagonal matrix  $\text{diag}\{\sin(y_i), i = 1, \dots, d\}$ . Any vector  $y$  such that, for all  $i = 1, \dots, d$ ,  $y_i = k_i \pi$  with  $k_i \in \mathbb{Z}$  is thus a stationary solution, *i.e.*, the reparameterization renders any point on the boundary of  $\mathcal{X}$  stationary. However, a local optimization algorithm applied to  $\tilde{h}(\cdot)$  provides a local maximizer for  $h(\cdot)$ . Indeed, suppose that the maximization of  $\tilde{h}(\cdot)$  with respect to  $y \in \mathbb{R}^d$  yields a vector  $y^*$  such that

$$\sin(y_i^*) \neq 0 \text{ for } i = 1, \dots, q \text{ and } \sin(y_i^*) = 0 \text{ for } i = q + 1, \dots, d.$$

Since  $\tilde{h}(\cdot)$  is maximized, one may assume that

$$\frac{\partial^2 \tilde{h}(y)}{\partial y \partial y^\top} \Big|_{y=y^*} \text{ is semi-negative definite}$$

and

$$(\hat{y}(i) - y^*)^\top \frac{\partial \tilde{h}(y)}{\partial y} \Big|_{y=\hat{y}(i)} < 0, \quad \hat{y}(i) = y^* + \alpha_i e_i, \quad i = 1, \dots, d, \quad (15)$$

with  $e_i$  the  $i$ -th basis vector, for  $|\alpha_i|$  small enough. Denote  $x^* = \cos(y^*)$  (notice that  $x^* \in \mathcal{X}$ ) and take any  $i \in \{q+1, \dots, d\}$ ; (15) implies that  $\partial h(x)/\partial x_i|_{x=x^*} > 0$  when  $y_i^* = 2k\pi$  (i.e.,  $x_i^* = 1$ ) and  $\partial h(x)/\partial x_i|_{x=x^*} < 0$  when  $y_i^* = (2k+1)\pi$  (i.e.,  $x_i^* = -1$ ),  $k \in \mathbb{Z}$ . Therefore, there exist  $d$  coefficients  $\lambda_i \geq 0$  and  $\lambda'_i \geq 0$  such that the Lagrangian

$$\mathcal{L}(x, \lambda, \lambda') = h(x) + \sum_{i=1}^d \lambda_i (1 - x_i) + \sum_{i=1}^d \lambda'_i (1 + x_i)$$

satisfies  $\partial \mathcal{L}(x, \lambda, \lambda')/\partial x|_{x=x^*} = 0$  with  $\lambda_i(1 - x_i) = \lambda'_i(1 + x_i) = 0$  for all  $i$  (take  $\lambda_i = \lambda'_i = 0$  for  $i = 1, \dots, q$ ;  $\lambda_i = \partial h(x)/\partial x_i|_{x=x^*}$ ,  $\lambda'_i = 0$ , when  $x_i^* = 1$  and  $\lambda_i = 0$ ,  $\lambda'_i = -\partial h(x)/\partial x_i|_{x=x^*}$ , when  $x_i^* = -1$  for  $i = q+1, \dots, d$ ). Moreover, for any  $z \in \mathbb{R}^d$  such that  $z_i = 0$  for  $i = q+1, \dots, d$ ,

$$z^\top \frac{\partial^2 \mathcal{L}(x, \lambda, \lambda')}{\partial x \partial x^\top} \Big|_{x=x^*} z = z^\top D_{1\dots q}^{-1}(y^*) \left\{ \frac{\partial^2 \tilde{h}(y)}{\partial y \partial y^\top} \Big|_{y=y^*} \right\}_{1\dots q, 1\dots q} D_{1\dots q}^{-1}(y^*) z \leq 0.$$

The inequality is strict when the matrix given by the first  $q$  rows and  $q$  columns of the Hessian matrix  $\partial^2 \tilde{h}(y)/(\partial y \partial y^\top)|_{y=y^*}$  is negative definite, which is a sufficient condition for the (local) optimality of  $x^*$ , see, e.g., Luenberger (1973, p. 235).

## Appendix C: measures of performance

In situations where the optimal design is known, to measure the performance of our algorithm we can naturally compare the value of the design criterion  $\phi_{k_{\max}} = \phi(\xi_{k_{\max}})$ , with  $\xi_{k_{\max}}$  the design measure obtained when the algorithm stops, to the optimal value  $\phi^* = \max_{\xi \in \Xi(\mathcal{X})} \phi(\xi)$ . When there exists a unique optimal measure  $\xi^* \in \Xi(\mathcal{X})$  such that  $\phi(\xi^*) = \phi^*$ , it is also instructive to evaluate how close  $\xi_{k_{\max}}$  is to  $\xi^*$ .

For the results of Sect. 4 we use the Wasserstein metric  $W_p(\xi_{k_{\max}}, \xi^*)$  with  $p = 1$  and  $p = 2$  and the LévyProkhorov metric  $L(\xi_{k_{\max}}, \xi^*)$ ; both define metrizations for the topology of weak convergence.

For  $\nu_0$  and  $\nu_1$  two probability measures on  $\mathcal{X}$ , the Wasserstein metric  $W_p(\nu_0, \nu_1)$  is defined by

$$W_p(\nu_0, \nu_1) = \left( \inf_{\omega \in \Omega(\nu_0, \nu_1)} \int_{\mathcal{X} \times \mathcal{X}} \|x - x'\|^p d\omega(x, x') \right)^{1/p}, \quad p \geq 1,$$

with  $\Omega(\nu_0, \nu_1)$  the set of all couplings of  $\nu_0$  and  $\nu_1$ , *i.e.*, of all measures  $\omega$  on  $\mathcal{X} \times \mathcal{X}$  having  $\nu_0$  and  $\nu_1$  as marginals. When  $d = 1$  we have

$$W_p(\nu_0, \nu_1) = \left( \int_0^1 [\mathbb{F}_0^{-1}(t) - \mathbb{F}_1^{-1}(t)]^p dt \right)^{1/p}, \quad (16)$$

where  $\mathbb{F}_i^{-1}$  denotes the quantile function  $\mathbb{F}_i^{-1}(t) = \inf\{x : \mathbb{F}_i(x) \geq t\}$  with  $\mathbb{F}_i$  the distribution function for  $\nu_i$ ,  $i = 0, 1$ , see Major (1978, Theorem 8.1).

The Lévy-Prokhorov metric  $L(\nu_0, \nu_1)$  is defined by

$$L(\nu_0, \nu_1) = \inf\{\epsilon > 0 : \nu_0(\mathcal{A}) \leq \nu_1(\mathcal{A}^\epsilon) + \epsilon \text{ and } \nu_1(\mathcal{A}) \leq \nu_0(\mathcal{A}^\epsilon) + \epsilon \text{ for all sets } \mathcal{A} \subset \mathcal{X}\}$$

where  $\mathcal{A}^\epsilon = \{x \in \mathcal{X} : \text{there exists } x' \in \mathcal{A}, \|x' - x\| < \epsilon\}$ . We always use the  $L_2$ -distance in  $\mathcal{X}$ , *i.e.*,  $\|x - x'\| = [\sum_{i=1}^d (\{x\}_i - \{x'\}_i)^2]^{1/2}$ . Note that  $W_p(\nu_0, \nu_1)$  and  $L(\nu_0, \nu_1)$  are sensitive to the scaling of  $\mathcal{X}$ ; in particular,  $W_p(\nu'_0, \nu'_1) = \alpha W_p(\nu_0, \nu_1)$  when  $\nu'_0$  and  $\nu'_1$  are obtained from  $\nu_0$  and  $\nu_1$  by applying the homothetic transformation  $x \in \mathcal{X} \mapsto z = \alpha x$ ,  $\alpha > 0$ .

Both metrics are difficult to compute for general probability measures but the fact that the measures  $\xi^*$  and  $\xi_{k_{\max}}$  have finite support makes the situation much simpler. Computational details will be presented elsewhere.

## Acknowledgments

The authors thank the referees for useful comments that helped to significantly improve the paper.

## References

- Atwood, C., 1973. Sequences converging to  $D$ -optimal designs of experiments. *Annals of Statistics* 1 (2), 342–352.
- Atwood, C., 1976. Convergent design sequences for sufficiently regular optimality criteria. *Annals of Statistics* 4 (6), 1124–1138.
- Ben-Tal, A., Nemirovski, A., 2001. *Lectures on Modern Convex Optimization: Analysis, Algorithms, and Engineering Applications*. MPS/SIAM Series on Optimization 2, Philadelphia.
- Böhning, D., 1985. Numerical estimation of a probability measure. *Journal of Statistical Planning and Inference* 11, 57–69.
- Böhning, D., 1986. A vertex-exchange-method in  $D$ -optimal design theory. *Metrika* 33, 337–347.

- Boyd, S., Vandenberghe, L., 2004. *Convex Optimization*. Cambridge University Press, Cambridge.
- den Hertog, D., 1994. *Interior Point Approach to Linear, Quadratic and Convex Programming*. Kluwer, Dordrecht.
- Dette, H., Melas, V., 2011. A note on de la Garza phenomenon for locally optimal designs. *Annals of Statistics* 39 (2), 1266–1281.
- Dette, H., Pepelyshev, A., Zhigljavsky, A., 2008. Improving updating rules in multiplicative algorithms for computing  $D$ -optimal designs. *Journal of Statistical Planning and Inference* 53, 312–320.
- Dette, H., Studden, W., 1993. Geometry of  $E$ -optimality. *Annals of Statistics* 21 (1), 416–433.
- Fedorov, V., 1972. *Theory of Optimal Experiments*. Academic Press, New York.
- Fellman, J., 1989. An empirical study of a class of iterative searches for optimal designs. *Journal of Statistical Planning and Inference* 21, 85–92.
- Frank, M., Wolfe, P., 1956. An algorithm for quadratic programming. *Naval Res. Logist. Quart.* 3, 95–110.
- Gribik, P., Kortanek, K., 1977. Equivalence theorems and cutting plane algorithms for a class of experimental design problems. *SIAM J. Appl. Math.* 32 (1), 232–259.
- Harman, R., Pronzato, L., 2007. Improvements on removing non-optimal support points in  $D$ -optimum design algorithms. *Statistics & Probability Letters* 77, 90–94.
- Harman, R., Trnovská, M., 2009. Approximate  $D$ -optimal designs of experiments on the convex hull of a finite set of information matrices. *Mathematic Slovaca* 59 (5), 693–704.
- Harville, D., 1997. *Matrix Algebra from a Statistician's Perspective*. Springer, Heidelberg.
- Hiriart-Urruty, J., Lemaréchal, C., 1993. *Convex Analysis and Minimization Algorithms*, part 1 and 2. Springer, Berlin.
- Karlin, S., Studden, W., 1966a. Optimal experimental designs. *Annals of Math. Stat.* 37, 783–815.
- Karlin, S., Studden, W., 1966b. *Tchebycheff Systems: With Applications in Analysis and Statistics*. Wiley, New York.

- Kelley, J., 1960. The cutting plane method for solving convex programs. *SIAM Journal* 8, 703–712.
- Kiefer, J., Wolfowitz, J., 1960. The equivalence of two extremum problems. *Canadian Journal of Mathematics* 12, 363–366.
- Luenberger, D., 1973. *Introduction to Linear and Nonlinear Programming*. Addison-Wesley, Reading, Massachusetts.
- Major, P., 1978. On the invariance principle for sums of independent identically distributed random variables. *J. Multivariate Analysis* 8, 487–517.
- McCormick, G., Tapia, R., 1972. The gradient projection method under mild differentiability conditions. *SIAM Journal of Control* 10 (1), 93–98.
- Molchanov, I., Zuyev, S., 2001. Variational calculus in the space of measures and optimal design. In: Atkinson, A., Bogacka, B., Zhigljavsky, A. (Eds.), *Optimum Design 2000*. Kluwer, Dordrecht, Ch. 8, pp. 79–90.
- Molchanov, I., Zuyev, S., 2002. Steepest descent algorithm in a space of measures. *Statistics and Computing* 12, 115–123.
- Nesterov, Y., 2004. *Introductory Lectures to Convex Optimization: A Basic Course*. Kluwer, Dordrecht.
- Nesterov, Y., Nemirovskii, A., 1994. *Interior-Point Polynomial Algorithms in Convex Programming*. SIAM, Philadelphia.
- Powell, M., 1964. An efficient method for finding the minimum of a function of several variables without calculating derivatives. *Computer J.* 7, 155–162.
- Pronzato, L., 2013. A delimitation of the support of optimal designs for Kiefer’s  $\phi_p$ -class of criteria. *Statistics & Probability Letters* 83, 2721–2728.
- Pronzato, L., Pázman, A., 2013. *Design of Experiments in Nonlinear Models. Asymptotic Normality, Optimality Criteria and Small-Sample Properties*. Springer, LNS 212, New York, Heidelberg.
- Pukelsheim, F., 1993. *Optimal Experimental Design*. Wiley, New York.
- Sadagopan, S., Ravindran, A., 1982. A vertex ranking algorithm for the fixed-charge transportation problem. *Journal of Optimization Theory and Applications* 37 (2), 221–230.

- Schwabe, R., 1995. Designing experiments for additive nonlinear models. In: Kitsos, C., Müller, W. (Eds.), MODA4 – Advances in Model-Oriented Data Analysis, Spetses (Greece), june 1995. Physica Verlag, Heidelberg, pp. 77–85.
- Schwabe, R., 1996. Optimum Designs for Multi-Factor Models. Springer, New York.
- Sergeyev, Y., Kvasov, D., 2006. Global search based on efficient diagonal partitions and a set of Lipschitz constants. SIAM J. Optim. 16 (3), 910–937.
- Sibson, R., Kenny, A., 1975. Coefficients in  $D$ -optimal experimental design. Journal of Royal Statistical Society B37, 288–292.
- Silvey, S., 1980. Optimal Design. Chapman & Hall, London.
- Silvey, S., Titterton, D., Torsney, B., 1978. An algorithm for optimal designs on a finite design space. Commun. Statist.-Theor. Meth. A7 (14), 1379–1389.
- St. John, R., Draper, N., 1975.  $D$ -optimality for regression designs: a review. Technometrics 17 (1), 15–23.
- Titterton, D., 1976. Algorithms for computing  $D$ -optimal designs on a finite design space. In: Proc. of the 1976 Conference on Information Science and Systems. Dept. of Electronic Engineering, John Hopkins University, Baltimore, pp. 213–216.
- Torsney, B., 1983. A moment inequality and monotonicity of an algorithm. In: Kortanek, K., Fiacco, A. (Eds.), Proc. Int. Symp. on Semi-infinite Programming and Applications. Springer, Heidelberg, pp. 249–260.
- Torsney, B., 2009. W-iterations and ripples therefrom. In: Pronzato, L., Zhigljavsky, A. (Eds.), Optimal Design and Related Areas in Optimization and Statistics. Springer, Ch. 1, pp. 1–12.
- Wu, C., 1978a. Some algorithmic aspects of the theory of optimal designs. Annals of Statistics 6 (6), 1286–1301.
- Wu, C., 1978b. Some iterative procedures for generating nonsingular optimal designs. Comm. Statist. Theory and Methods A7 (14), 1399–1412.
- Wynn, H., 1970. The sequential generation of  $D$ -optimum experimental designs. Annals of Math. Stat. 41, 1655–1664.
- Yang, M., 2010. On de la Garza phenomenon. Annals of Statistics 38 (4), 2499–2524.
- Yang, M., Stufken, J., 2009. Support points of locally optimal designs for nonlinear models with two parameters. Annals of Statistics 37 (1), 518–541.

- Yu, Y., 2010a. Monotonic convergence of a general algorithm for computing optimal designs. *Annals of Statistics* 38 (3), 1593–1606.
- Yu, Y., 2010b. Strict monotonicity and convergence rate of Titterington’s algorithm for computing  $D$ -optimal designs. *Comput. Statist. Data Anal.* 54, 1419–1425.
- Yu, Y., 2011.  $D$ -optimal designs via a cocktail algorithm. *Stat. Comput.* 21, 475–481.